

Gene-Environment Interaction Analysis using GPU

Daniel Berglund

November 2014

Abstract

Abstract

Contents

1	Introduction	1
1.1	Outline	1
1.2	Genome-wide association studies	1
1.3	Genetics	2
1.3.1	Major, Minor and Risk Allele	2
1.4	Defining Interaction	3
1.5	GEIRA, JEIRA and GEISA, the Aim of the Project	4
2	Background	5
2.1	Statistical Background	5
2.1.1	Contingency Tables	5
2.1.2	Relative Risk and Odds Ratio	5
2.1.3	Additive Interaction	6
2.1.4	Statistic Measures for Additive Interaction	6
2.1.5	Confounders and Covariates	8
2.1.6	The Multiple Testing Problem	8
2.2	Computer Architecture Background	9
2.2.1	Central Processing Unit	9
2.2.2	Concurrency and Threads	9
2.2.3	Memory, Caching and Optimizations	11
2.2.4	Accelerators, GPU and Xeon Phi	15
2.2.5	Data Storage	17
2.2.6	Clusters	17
2.3	CUDA programming model	20
2.3.1	Device Memory	21
2.3.2	Streams	23
2.3.3	Efficient CUDA	24
2.4	Software Design Background	27
2.4.1	Unit Tests and Mocks	28
2.4.2	Design Patterns	28
2.4.3	Version Control Software	29
2.5	Performance Measures	30
2.5.1	Amdahl's Law and Gustafson's Law	31
2.5.2	Profilers	33

3 Algorithms and Implementation	34
3.1 Algorithmic Background	34
3.1.1 Logistic Regression	35
3.1.2 Data Mining and Machine Learning Approaches	36
3.2 Implementation	39
3.2.1 JEIRA and GEISA	39
3.2.2 CuEira	39
3.2.3 Kernels	44
3.2.4 GPU stuff	44
4 Results	46
5 Discussion and Conclusions	47
6 Outlook	48
Bibliography	49
Appendices	54
A Lists	55
List of Algorithms	57
List of Examples	58
List of Figures	59
List of Tables	60
B File Formats	61
B.1 Plink Binary Data Format	61
B.2 Environmental and Covariates File Format	61
C Bugs and Issues	62
C.1 Risk Allele Code from GEIRA, JEIRA and GEISA	62
C.2 Encountered Bugs	62
D How to Compile and Use CuEira	64

Chapter 1

Introduction

1.1 Outline

The first chapter contains a short introduction to the area and the terminology. The second chapter gives the theoretical background. It starts with the statistics and algorithms and then describes the computer architecture and software design. The third algorithm section explains the structure of the current algorithms and also explains the design of the program that was developed in this thesis called CuEira. Section four provides the results and comparisons of the program performance. The fifth section is the discussion of these results. The last section discusses the future work.

1.2 Genome-wide association studies

One type of study to find associations between genetic markers and diseases or other traits in different individuals is genome-wide association study(GWAS). Most GWAS do not study interactions between the genetic markers or between the genetic markers and environmental factors [1, 2]. Investigating gene-gene interactions recently has become more common[1], however gene-environment interactions are still uncommon in studies[2]. Interactions between genes and environmental factors are considered to be important for complex diseases such as cancer or autoimmune diseases[1, 2, 3, 4]. In general a complex disease develops due to a combination of factors, not just a single gene or environmental factor[5]. Examples of environmental factors are smoking, physical activity and so on.

GWAS usually has a study design that is either cohort or case-control[5]. In cohort study a sample of a population is followed and over time some individuals will develop the disease of interest[6]. In case-control studies two groups are compared with each other to find the risk factors[6]. One of these two groups consists of individuals with the disease and the other group consists of individuals similar to the first group but who do not have the disease[6].

A typical GWAS consists of tens of thousands of individuals and up to several millions of genetic markers[1, 7]. In gene-gene interaction GWAS few studies investigate more than the second order interaction because of the high number

of possible combinations. There are some algorithms that can investigate higher order interactions however these have drawbacks[8, 9, 10].

1.3 Genetics

The genetic information is stored in *chromosomes* each consisting of a *DNA molecule*[11]. Each chromosome comes in a pair where both chromosomes are nearly identical, except for the chromosomes related to gender[11]. The DNA molecule is a double stranded helix of four nucleotide bases, *adenine(A)*, *cytosine(C)*, *guanine(G)* and *thymine(T)*. They are always paired as A-T and G-C. *Genes* are sections of these DNA molecules. The pair of genes in the chromosomes is a *genotype*. The observable product of the genotype is called *phenotype*, e.g. eye colour, fur patterns.

A position of the DNA or a gene is a *locus*[11]. A variation of the same gene or locus is an *allele*[11]. The effect of an allele on a phenotype can be either *dominant*, *recessive* or *co-dominant*. Dominant means that if the allele is present the phenotype is expressed. For recessive the allele needs to be present in both chromosomes to be expressed[11]. Co-dominant is when both alleles are expressed, when different alleles are present this usually produce some kind of intermediate state[11]. An example of co-dominance is the human blood groups[11].

Genotype	Dominant	Recessive
AA	Effect	Effect
GA	Effect	No Effect
GG	No effect	No effect

Table 1.1: The phenotype based on the genetic model and if the allele with the effect is A

1.3.1 Major, Minor and Risk Allele

The the allele that is least common using all individuals(both cases and controls) is the *minor allele*, the one that is most common is the *major allele*. The frequency of the minor allele is the *minor allele frequency(MAF)*[12, 3]. If the MAF is too low an analysis often will be of little value so all SNPs with MAF below a threshold are usually skipped. 5% is common value of the threshold[7, 3].

To determine if the genetic risk is present for each individual one of the alleles needs to be chosen as the *risk allele*. GEIRA and JEIRA determines the risk allele by comparing the MAFs of case and control[13, 12]. If MAF of cases is greater or equal to the MAF of controls the minor allele is used as the risk otherwise the other allele[13, 12].

A simpler equivalent definition is that the risk allele is the allele which frequency is higher in case than in controls. This definition means that the allele chosen as risk might not be the most common in either case or control, however the portion of the population that has the allele is higher among cases than controls. This hints that it might be a reason behind the outcomes. This

definition is equivalent to the definition using MAF because if the minor allele has a higher frequency in cases than it the risk is the minor allele according to both definitions. If that is not the case the first definition sets the major allele as the risk. Also due to symmetrical reasons(the total frequency is always 1) this means that the major allele has higher frequency in cases than in controls, which according to the second definition means it is the risk allele.

The presence of genetic risk based on the risk allele and the genetic model can be coded as a variable than then be can used in a statistical model. For dominant or recessive genetic model this means that the variable is binary since the risk is either present or not. See table 1.1 for how the coding is based on the model. There is an alternative way to code using co-dominant genetic model. In this case the number of risk alleles in the individual used as the variable, so either 0, 1 or 2.

1.4 Defining Interaction

There are several ways to define interaction. The overall goal is often to detect if *biological* interactions are present. *Biological* interaction is when different factors co-operate through a physiological or biological mechanism and cause the effect, e.g. the disease. The information about the factors can then be used to explain the mechanisms involved in causing the disease and possibly help to find cures for them. However biological interaction is not well defined and thus it is not possible to calculate it directly from data.[14, 5, 15]

That is why statistical interaction is used and it is assumed that the presence of statistical interaction implies the presence of biological interaction

Statistical interaction on the other hand is well defined. However it is scale dependent, i.e. interactions can appear and disappear based on transformations of the given data. Statistical interaction also depends on the model used. The common way to define statistical interaction is to consider the presence of a product term between the variables in the model, this is referred to as *multiplicative* interaction. For instance for a linear model

$$f(x, y) = ax + by + cxy + d \quad (1.1)$$

c is the product term that represents multiplicative interaction between variables x and y so the test for multiplicative interaction is to test if $c=0$ [3, 14, 5]. In some cases it is instead tested as $a=b=c=0$, this is testing for *association allowing for interaction*[1].

Additive interaction is another kind of statistical interaction that is broader than multiplicative interaction. It is when the total effect of the interacting factors present is larger than the sum of their separate effects. A more precise definition can be found in section 2.1.3. It implies biological interaction as defined by Rothman[5], which is sometimes called *causal interdependence* or *causal interaction*[16]. Additive interaction can find more possible interactions than multiplicative and since it is closer to biological interaction it is often chosen in epidemiological studies[14].

1.5 GEIRA, JEIRA and GEISA, the Aim of the Project

To search for interaction in data various tools can be used. *GEIRA*[3] is one such tool which uses logistic regression, additive interaction and with bootstrap to asses the quality of the models[3]. It comes in two versions, one made in R and one in SAS. To improve speed up it was later implemented in Java and called *JEIRA*[12]. It is parallel which made it significantly faster[12]. *GEISA*[13] is another tool based on JEIRA however it uses permutation testing instead of bootstrap to asses the certainty of the models[13]. The programs structure is explained in section 3.2, the algorithms are explained in section 3.1. Bootstrap calculates the certainty by resampling the data[17]. Permutation tests works in a similar way but only randomizes the outcomes [18].

The aim of this project is to make the gene-environment interaction analysis faster by using graphics processors and to be able to handle larger amounts of data. The program is written in C++ and CUDA.

Chapter 2

Background

This section will go through some of it starting with the algorithms and statistical background. And continuing with computer architecture and software design after that.

2.1 Statistical Background

This section will explain some of the statistical background of categorical statistics.

2.1.1 Contingency Tables

A contingency table is a matrix used to describe categorical data. Each cell contains a count of occurrences for a specific combination of variables. Table 2.1 is an example of an 2×2 table. From this table we can for instance see that 688 smokers got lung cancer. Contingency tables are the basis for various statistical tests to model the data.[17]

	Lung cancer	No lung cancer
Smoker	688	650
Non smoker	21	59

Table 2.1: Contingency table describing the outcome of a study, from [17], page 42

2.1.2 Relative Risk and Odds Ratio

From a contingency table it is possible to calculate some useful measures such as the risk of getting the outcome based on exposure. The risk for a row i in the table will be referred to as π_i . For table 2.1 the risks π_1 and π_2 is

$$\pi_1 = \frac{688}{688 + 650} = 0.51 \quad (2.1)$$

$$\pi_2 = \frac{21}{21 + 59} = 0.36 \quad (2.2)$$

To compare different risks, for instance between smokers and non smokers, the ratio of the risks is used[17]. It is called *relative risk*(RR) is defined as[17]

$$RR = \frac{\pi_1}{\pi_2} \quad (2.3)$$

So for the table 2.1 the relative risk of getting lung cancer based on exposure to smoking is

$$RR = \frac{0.52}{0.36} = 1.96 \quad (2.4)$$

This means that the risk of getting lung cancer for a smoker is almost twice as high for a non smoker in this data.

Another useful measure is *odds* and *odds ratio*(OR). The odds is[17]

$$\Omega = \frac{\pi}{1 - \pi} \quad (2.5)$$

The odds are non-negative and $\Omega > 1$ when the outcome is more likely than not[17]. So for $\pi = 0.75$ the odds is $\Omega = \frac{0.75}{1-0.75} = 3$. This means that the outcome is 3 times more likely to occur than not. ORs can be used as an approximation of RR in case control studies because RR can not be estimated in that type of studies[19]. Cohort studies on the other hand can give estimates of RR[19]. OR is the ratio of the odds just as RR is the ratio of the risks[17].

$$\theta = \frac{\Omega_1}{\Omega_2} \quad (2.6)$$

In the case when the outcome is a disease or similar variables with odds ratio below one are called *protective* and when it is above one it is a *risk factor*[20].

2.1.3 Additive Interaction

Interaction was talked about in section 1.4. The more precise definition of additive interaction is the divergence from additive effects on a logarithmic scale, e.g[5].

$$OR_{both\ factors\ present} > OR_{first\ factor\ present} + OR_{second\ factor\ present} - 1 \quad (2.7)$$

A third variable is used to model the interaction for both multiplicative and additive interaction. Multiplicative interaction as explained in section 1.4 does not include the effects from the factors, i.e. the main effects, in the interaction variable. However for additive interaction the interaction variable models the whole effect including the main effects. This means that the third variable for multiplicative and additive interaction are identical, however for additive interaction the first and second variable is zero when both factors are present. The coding for additive interaction is summarized in table 2.2[12].

2.1.4 Statistic Measures for Additive Interaction

Based on the statistical model and its corresponding ORs there are some measures that can be calculated. These measures show various properties of the additive interaction. They are defined using *RR* however as mentioned before

Factor present	First variable	Second variable	Interaction
None	0	0	0
First	1	0	0
Second	0	1	0
Both	0	0	1

Table 2.2: Coding of the variables for LR

in section 2.1.2 OR can be used to approximate RR .

Relative excess risk due to interaction(RERI) is how much of the risk is due to interaction[5, 20]. It is defined as[5, 20]

$$RERI = RR_{11} - RR_{10} - RR_{01} + 1 \quad (2.8)$$

Attributable proportion due to interaction(AP) is similar to RERI however is the proportion relative to the interaction relative risk[5, 20].

$$AP = \frac{RERI}{RR_{11}} = \frac{1}{RR_{11}} - \frac{RR_{10}}{RR_{11}} - \frac{RR_{01}}{RR_{11}} + 1 \quad (2.9)$$

Synergy index(SI) is the ratio of the combined effects and the individual effects[5, 20].

$$\frac{RR_{11} - 1}{RR_{10} + RR_{01} - 2} \quad (2.10)$$

By using the definition of additive interaction it can be shown that additive interaction is present when $RERI > 0$, $AP > 0$ and $SI > 0$ [20].

2.1.4.1 Recoding of Protective Effects

The measures for additive interaction RERI, AP and SI explained in the previous section are developed for risk factors, i.e. $OR > 1$, which causes problems if a factor is protective[20]. This can be solved by *reencoding*, it switches the group of individuals without both factors, i.e. the *reference group*, with the group that has the lowest OR [20]. There are three possible combinations of the factors excluding the reference group shown in table 2.3. Switching a group changes the level of the factors for all in individuals in the group to the reference groups level, and switches the individuals in the reference group with the groups level.

Genetic factor	Environment factor	Reencoding case
0	0	N/A
1	0	1
0	1	2
1	1	3

Table 2.3: Cases of recoding

2.1.5 Confounders and Covariates

Confounding is one of the central issues in design of epidemiological studies. It is when the effect of the exposure is mixed with the effect of another variable. So if we do not measure the second variable, the effect of the first would be estimated as stronger than it really is. The second variable is then a *confounder*. Several methods in epidemiology are about avoiding or adjusting for confounding. Sometimes these variables need to be incorporated into the models. *Covariates* are possible confounders or other variables that one wants to adjust for in the model. Sometimes covariates are called *control variables*.[15, 5]

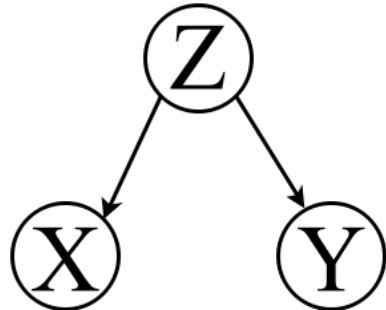


Figure 2.1: Illustration of a simple case of confounding. If we do not observe Z we might falsely find an association between X and Y

For instance if a study would look into the effects of yellow teeth on lung cancer they would likely find an effect. However it would be due to the confounder smoking since smoking causes yellow teeth and lung cancer.

2.1.6 The Multiple Testing Problem

It is not uncommon to test many hypothesis on the same data and in the case of GWAS it can be billions of tests. If one continues to test one should eventually find something that is significant due to random chance. With the common significance threshold of 5% it is expected to get 1 in 20 false positives under the assumption that the null hypothesis is true. The problem arises from the fact that the hypothesis tests are dependent on each other since they use parts of the same data. This is the multiple testing problem and if it is not corrected for many false positives might be found.[21]

Bonferroni correction is the simplest method and viewed as a conservative way to correct for this problem. It simply divides the significance threshold by the number of hypothesis tested. However the number of hypothesis made is not always clear. With a two-stage analysis, is the number of hypothesis the number of tests done in both stages combined, the number made in the first stage or the second stage?[21]

2.2 Computer Architecture Background

This section explains some of the computer architectures used to perform tasks in a thesis. A large part of this chapter is focused on the optimization techniques that are in use. For the purpose of this thesis it is best to think of the computer consisting of four main parts, the processor, the data storage, the memory and the accelerator.

2.2.1 Central Processing Unit

Central processing unit(CPU) is the part of the computer that executes instructions one by one[22]. Most modern CPUs have multi-core architecture meaning they consist of several processors[22]. Each core can perform tasks independent of each other.

The figures 2.2 and 2.5 shows how the CPU is divided into areas and that most of it is not used for calculations. A large part of the area is used for various optimizations.

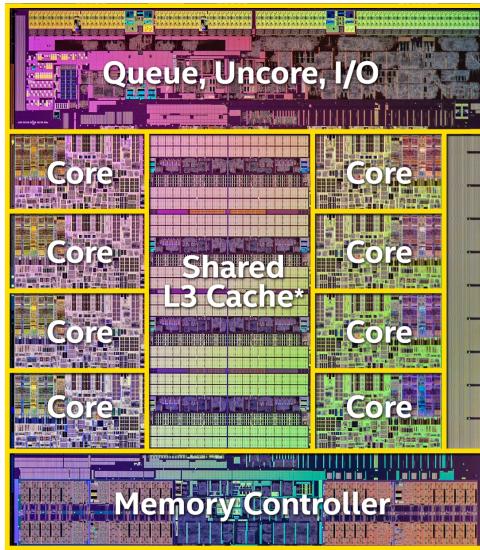


Figure 2.2: The layout of the Haswell architecture i7-5960x, from [23]

so the programs need to be parallel to get maximum speed. Parallel means that the program issues multiple instructions at the same time, more about this in section 2.2.2.

2.2.2 Concurrency and Threads

Concurrency means doing multiple things at the same time, instead of *sequential*. This can cause various problems such as the same object being accessed at the same time. The dinning philosophers is an concurrency problem that can be used to illustrate some of these problems. A group of philosophers is sitting at a round table, each has a plate of spaghetti in front of them and there is a

fork between each pair of philosophers[24]. The philosophers alternate between thinking and eating. However they can only eat if they have both the right and left fork[24].

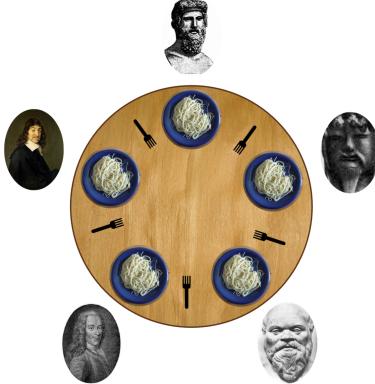


Figure 2.3: Illustration of the dining philosophers problem. Wikipedia Commons

A potential proposal[24] for behaviour instructions could be:

- Think
- Wait for a fork to become available and pick up that fork
- Wait for and pick up the other fork
- Eat
- Put down the forks one by one
- Go back to thinking

The problem is that this set of instructions can lead to a state where everyone is holding one fork and waiting to get the other one[24]. But no one is done eating so no forks will become available. The system is then locked into its state, this is called a *deadlock*[24, 22]. A potential solution to the problem in this case is to introduce a person that dictates if a philosopher is allowed to eat or not[24].

There are also other potential problems that can arise when using concurrency. A *race condition* occurs when the result depends on what affects it first, it is a race between them on who can reach it first[22]. *Locks* can be used to prevent situations as the ones described above by making sure only one thread at a time can access the object[22].

In computer science concurrency occurs when instructions are separated in different *threads*. The architectures can be divided in a taxonomy with four categories including the sequential. They are illustrated in figure 2.4.

- *SISD*(Single-Instruction,Single-Data) is the sequential and is as the name says, single instructions on single data[22].
- *SIMD*(Single-Instruction, Multiple-Data) is applying the same instruction on different data[22]. SIMD is also called vectorisation[22].

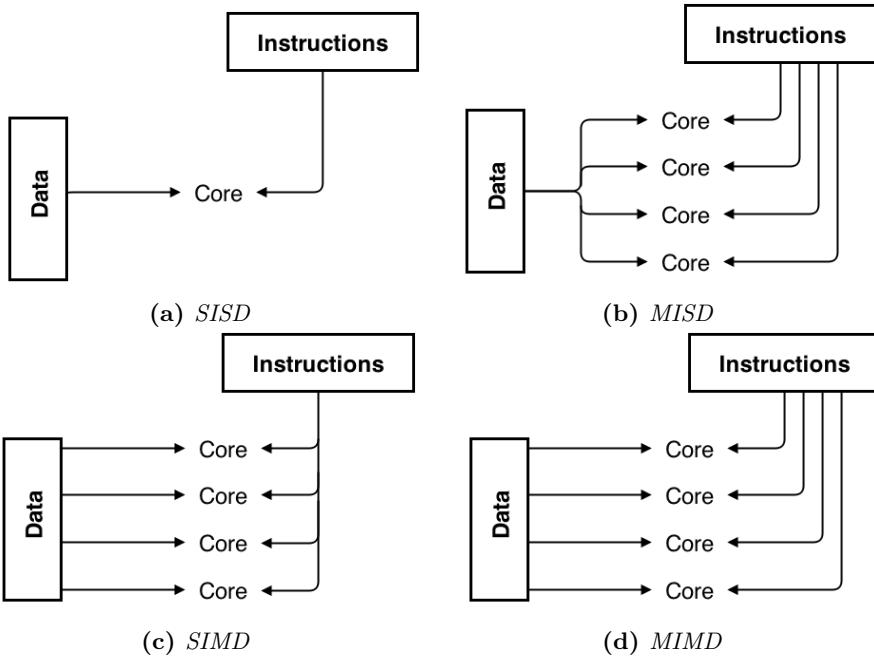


Figure 2.4: Flynn's taxonomy of parallel architectures[22]

- **MIMD**(Multiple-Instruction, Multiple-Data) applies different instructions on different data[22].
- **MISD**(Multiple-Instruction, Single-Data) performs different operations at the same piece of data, this paradigm is uncommon[22, 25].

2.2.3 Memory, Caching and Optimizations

Various optimizations have been introduced to the CPUs over the years[26, 22]. Some are common and almost all CPUs have them, others are unique to a specific vendor or CPU[22]. Figure 2.5 shows the layout of a CPU core. This section explains some of the more common optimizations.

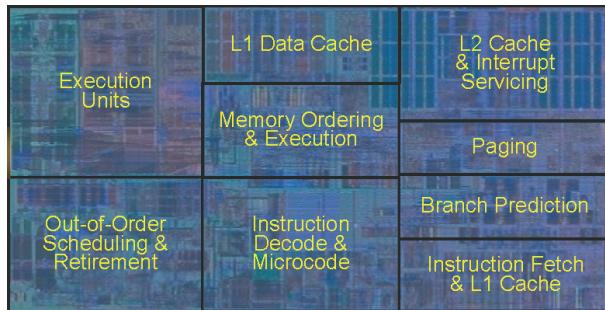


Figure 2.5: Layout of a CPU core, Intel Nehalem i7, from [27]

In general retrieving data from memory is relatively slow compared to how fast the processor works[22, 26]. To solve this problem a *cache* was introduced[26]. It stores recently used data in a very small and very fast memory that resides on the CPU chip itself[26]. The caches varies in sizes, the latest ones are around 16-20 mB, for instance Intels i7-5960x with 20 mB [28]. The cached data can then be reused without waiting for the main memory to fetch it. However when the data is not in the cache it costs time to fetch it from the main memory, that is a *cache miss*[26]. Avoiding cache misses is important for the speed up[26]. Modern CPUs commonly have more than one cache, most have three[22]. They are named L1 to L3, with L1 being the one fastest and smallest and L3 the largest and slowest[22].

Modern multi-core processors have several L1 caches, one per core[22]. Commonly each core has a L2 cache as well, however it can be shared in some architectures[22]. The L3 cache is almost always shared between all cores[22]. The multiple caches used causes a problem called *cache coherency*. It is when the data modified in one cache needs to be updated in all the caches in order for them to use the new value[22].

One of the memory optimization methods somewhat related to caches is *prefetching*[22, 26]. Instead of fetching just the data requested by the current operations it also fetches the surrounding data[26]. Even if the data is not needed at the given point in time, the chance that the surrounding data will be used soon is high since it is common to iterate over arrays and similar structures[22].

When a program iterates over an array the first value is a cache miss since it is not in the cache. The following values are prefetched and cached as well. However if the array is too long to be prefetched and cached completely and the CPU hits a position not in the cache then a new cache miss happens[26]. Usually matrices are stored sequential in memory[26]. For instance a two dimensional matrix is stored as a one dimensional array. Different programming languages store them differently, either column by column or row by row[26]. Row by row is called *row major* and column by column is *column major*.

Below is an example of how prefetching can have a large impact on performance.

In section 6.2.1 in What Every Programmer Should Know About Memory[26] an example of how much speed that can be gained by optimizing matrix matrix multiplication. The variables have been renamed for clarity. The simple version of it is shown in algorithm 1 and the version with improved cache use is shown in 2.

By transposing the mul2 matrix it becomes more cache friendly by using prefetching of surrounding data the way that mul1 already is[26]. The transpose itself can be eliminated while also making it use the cache better by reading in the correct amount of data, this is the size SM in the code[26]. The optimized version is shown in 2. It took 17.3% of the original time[26].

On top of that by using some additional features of the CPU such as doing multiple instructions at once more speed up is gained increasing it to 9.47% of the original time[26].

Algorithm 1: Basic algorithm for matrix matrix multiplication of two $N \times N$ matrices

Data: Three $N \times N$ matrices, $matrix1$, $matrix2$ and $result$

```
for (i = 0; i < N; ++i)
{
    for (j = 0; j < N; ++j)
    {
        for (k = 0; k < N; ++k)
        {
            result[i][j] += matrix1[i][k] * matrix2[k][j];
        }
    }
}
```

Algorithm 2: Matrix matrix multiplication optimized

Data: Three $N \times N$ matrices, $matrix1$, $matrix2$ and $result$

SM is the number of matrix elements that fits in the cache.

```
for (i = 0; i < N; i += SM)
{
    for (j = 0; j < N; j += SM)
    {
        for (k = 0; k < N; k += SM)
        {
            for (i2 = 0, rresult = result[i][j], rmatrix1 = matrix1[i][k];
                 i2 < SM; ++i2, rresult += N, rmatrix1 += N)
            {
                for (k2 = 0, rmatrix2 = matrix2[k][j]; k2 < SM; ++k2,
                     rmatrix2 += N)
                {
                    for (j2 = 0; j2 < SM; ++j2)
                    {
                        rresult[j2] += rmatrix1[k2] * rmatrix2[j2];
                    }
                }
            }
        }
    }
}
```

If the operation instructions diverge, e.g. the *if* statements, the CPU would have to wait for the results to see what instructions load and do next[26]. This can mean a significant loss of time, however because of an optimizations called *branch prediction* and *speculative execution* the loss of time can be prevented in some cases[26]. A dedicated part of the CPU stores the results from these divergences and when they are encountered again it uses them to make a guess of what do to next[26]. It then loads the instructions and tries to fetch the needed data[26]. The CPU then speculatively executes the instructions, i.e it executes instructions that might not actually be needed[26]. If the guess was correct the results are kept, on the other hand if it was wrong the CPU starts again from the correct path and discards the incorrect results[26].

Sometimes there can be several independent instructions, e.g. when adding two vectors together, this allows optimizations for *instruction level parallelism*[22]. The order of the instructions can be ignored which for instance means that if an instruction needs to wait for some data another instruction that has its data available could be executed while waiting[22]. It is an optimization called *out of order execution*[22]. There are also other possible optimizations[22].

In the example 1 the instructions are independent because they use different variables so they can be executed in any order. If the variables *b* and *c* needs to be fetched but *e* and *f* are available due to previous instructions then the CPU could execute the second addition first by using out of order execution. However if an operation later depends on the variable *a* or *d* then that instruction would have to wait until the additions have been executed.

Example 1: Independent Instructions

Data: Six int variables, *a*, .., *f*

$$\begin{aligned} a &= b + c \\ d &= e + f \end{aligned}$$

Pointer aliasing can prevent the compiler from making certain optimizations, e.g. out of order execution mentioned above. Pointer aliasing is when the same memory area can be accessed using different variables[22]. This becomes a problem if an set of instructions use these variables in the same area of the program. The relative order of operations between these operations then has to be preserved. The compiler does have the information if aliasing occurs or not which means the compiler has to treat all situations where it can occur as if it does[22].

If the integers in example 1 are changed to pointers then aliasing can occur because some of these pointers could point to the same memory. For instance if *a* and *e* points to the same memory then the result would be wrong if the second line is executed before the first. This would force the instructions to be executed in order or the result would likely be incorrect.

Some programming languages, e.g. Fortran, disallow some types of aliasing while others, e.g. C/C++, allows aliasing[22]. The compilers for the languages in the second case commonly have an option to disable aliasing. However for languages that don't allow aliasing, or if it is disabled, it is then the programmers

Example 2: Pointer Aliasing

Data: Six int pointers, a, \dots, f

$$\begin{aligned} *a &= *b + *c \\ *d &= *e + *f \end{aligned}$$

responsibility to make sure that aliasing does not occur or the results can be wrong[22].

2.2.4 Accelerators, GPU and Xeon Phi

Accelerators are separate parts made to do one task and do it well. They are usually highly specialised and perform badly outside of the tasks they were designed for. The two accelerators explained here is graphical processing unit(GPU) and Intel Xeon Phi[29, 30]. GPUs have mostly been used for games and other graphics related problems however have become more popular for general computing the last 10 years due to their high number of cores[29]. Intel Mic is much more recent with prototypes in 2010 and the first generation released in 2010[30]. It is Intels response to the GPUs and has elements of both GPU and CPU[30]. This section will focus on explaining some parts of the architecture of the GPUs and Xeon Phi. How to use the GPUs are talked about in more detail in section 2.3.

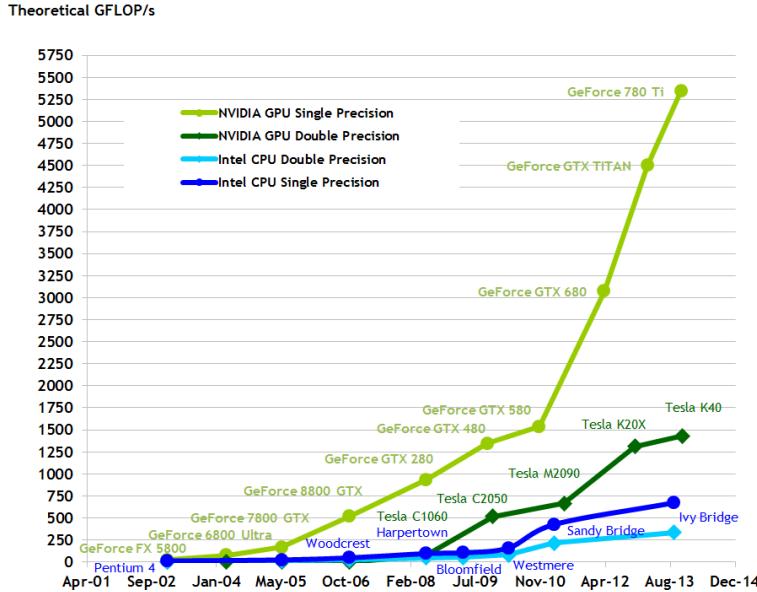


Figure 2.6: Theoretical throughput of some NVIDIA GPUs and Intel Processors, from[29]. However the chart does not contain a Xeon Phi processor.

GPUs are made to process images in a fast manner so that games and similar tasks runs smoothly[29]. Because GPUs are made for image processing they are GPUs have a large number of cores(1000-2000)[29, 32]. Each core is slower

than a core in a CPU, however the number of cores makes up for it. The downside is that it's slower than a normal CPU for sequential tasks so the program has to take advantage of the extreme parallelism to get good speed[33]. The memory of the GPU is separate from the computers main memory which means that the data needs to be transferred between the memories which increases the programs complexity[29].

The GPU cores work in a manner that is similar to SIMD[29]. The cores also have fewer optimizations than CPUs, this means that they can then devote a larger portion of the chip to the calculations as seen in figure 2.7. For instance NVIDIAs GPUs lack optimizations such as branch prediction and speculative execution talked about in section 2.2.3 [29]. Their caches are also much smaller than the CPUs[29]. Another inheritance from the image processing is that they are much faster with single precision than with double precision data types[29, 32]. A CPU drops slightly in performance when using double precision however much less than a GPU as seen in figure 2.6. Single and double precision refers to how many bytes are used to store the decimal numbers, fewer bytes means less numerical precision.

Intels answer to the increasing popularity of GPUs for general computing is Xeon Phi[30]. It has elements of GPU architecture however is still more like a CPU. It has separate memory and many cores(50+)[30]. It is MIMD which can be good for codes that diverge a lot, for ease of use however it has the possibly to use more SIMD like structures[30]. It is possible to only use the Xeon Phis memory and avoid the memory transfers that way. The memory is relatively small compared to the main memory so if the program needs a lot of memory transfers are likely needed. The cores of the Xeon Phi lacks some of the normal optimizations, e.g. out of order execution[30].

2.2.4.1 GPU vs Xeon Phi, Why use GPUs for Interaction?

Which is most efficient depends largely on the algorithm and if there is any previous existing code. For instance if the algorithm contains many points of divergence Xeon Phi is likely faster because Xeon Phi is MIMD while GPU is SIMT. Xeon Phi can be applied to a broader set of problems due to using a MIMD architecture[30]. A one on one comparison of speed might not be entirely relevant because the algorithm used for the test might simply be better suited for one the architectures. Speed per cost, either power consumption or the price on the hardware, can also be more relevant. If one is half the cost of the other then comparing two versus one would be a better comparison. CUDA is also more mature than Xeon Phi since Xeon Phi is only a few years old while CUDA is almost ten years old[29, 30].

Which suits best needs to be examined by a case by case basis and depending on the specifics of the algorithm and what type of computers that might already be available. It is also important to think about how important maximum speed is and how much time there is to spend. If extremely high speed is of importance and time is not a concern then a GPU is likely to get higher performance, unless the algorithm is badly suited for GPUs. However Xeon Phi might reach similar speed but with less effort.

GPUs are good for interaction algorithms because most methods consider each combination independently and also in this case that LR consists of linear algebra operations. A new combination can be transferred while another is calculated. Linear algebra with sufficiently large matrices and vectors work very well on GPUs[35, 29]. Several studies got high gains from implementing their algorithms on GPU[8, 36, 37, 38, 39, 40]. However most studies had CPU versions that were not parallel so it is likely that the gains would have been less if they would have been compared to an optimized and parallel CPU version, especially if a Xeon Phi co-processor was used. One study did make a comparison of a CPU cluster version and a GPU version of their algorithm, they found that 16 CPU nodes had the same performance as a single GTX 280 card[41].

2.2.5 Data Storage

2.2.6 Clusters

Clusters are collections of computers that can work together in several ways and are so tightly connected that they can usually be viewed as a single system. They communicate over a shared network but have separate memory and processors. Storage is usually shared. A computer in a cluster is called a *node*. [42, 22]

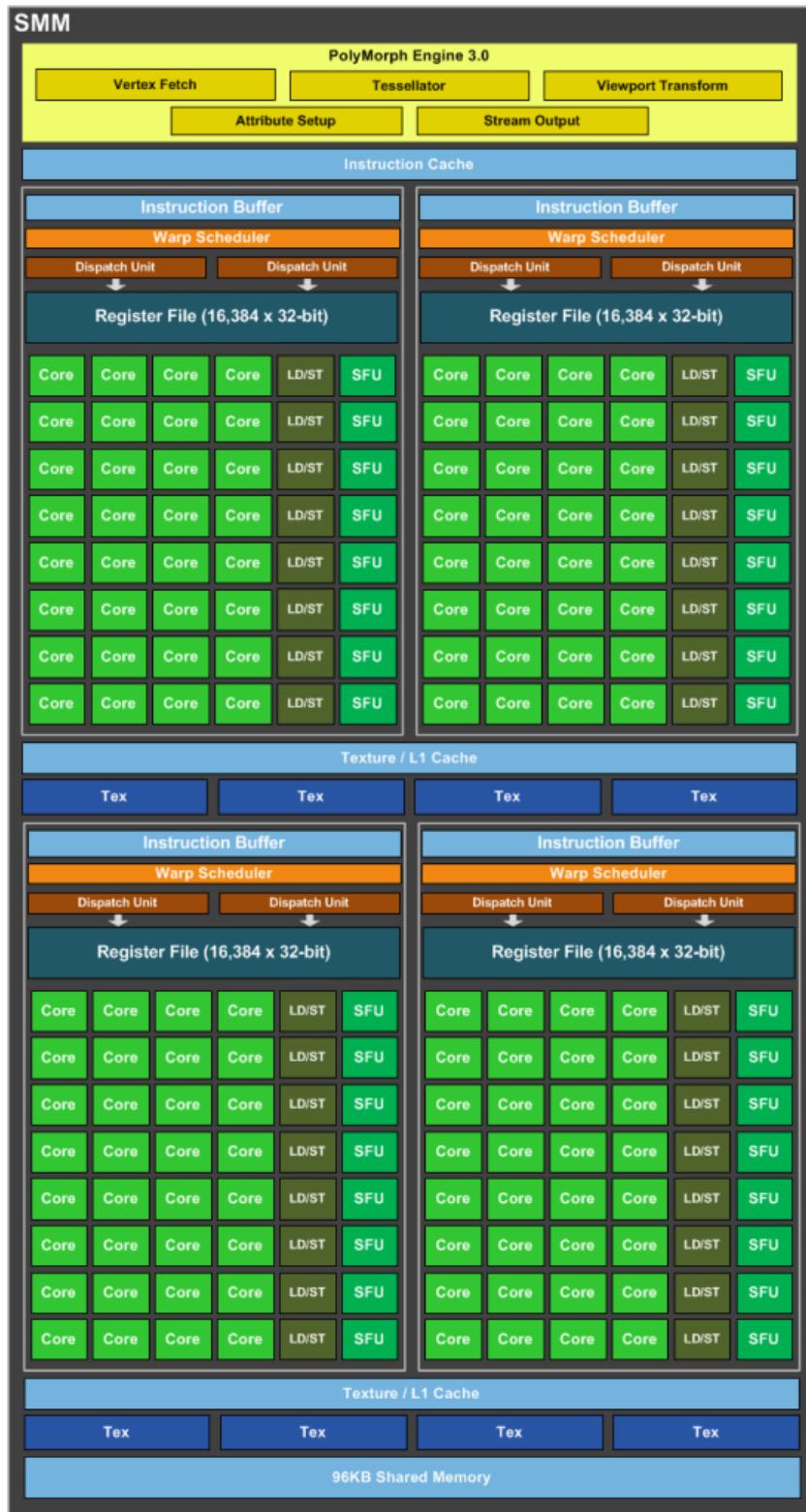


Figure 2.7: Layout of the Maxwell architecture, from [31]



Figure 2.8: Knights Corner silicon layout, from[34]

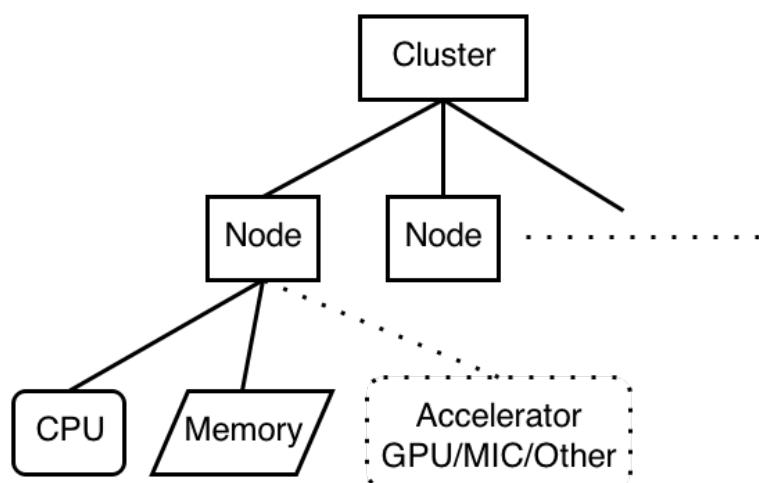


Figure 2.9: Basic overview of a cluster

2.3 CUDA programming model

This section is about GPU programming in more detail using NVIDIA's CUDA (Compute Unified Device Architecture)[29]. CUDA extends C and C++ with some additional functionality that can be used to perform operations on NVIDIA's GPUs. It provides a separate compiler to compile the GPU code called nvcc[29]. Different GPUs support different CUDA functions, each NVIDIA GPU has a value called computational capability. The higher the value the more features of CUDA are supported on that GPU. Some properties also vary between the different underlying architectures. This means that it is important to know what kind of capability the GPU to be used has so that the program does not need features that are not there. The things described here might not apply to GPUs with compute capability below 3. There are GPUs made specifically for calculations rather than games. This section is a summary of how CUDA works from the CUDA programming guide[29] and CUDA best practices guide[33].

In CUDA the GPU is called *device* and systems CPU and memory is the *host*. To perform operations on the device a type of function called *kernel* is used. A kernel works mostly as a normal C/C++ function with the addition of some specifiers to provide options to set number of threads and so on. Example 3 is an example of a simple kernel definition and example 4 shows how it can be called from the host code. The kernel adds each element of the arrays A and B storing it in C. Each thread performs one addition and knows which elements to use based on its id. The call to the kernel is made by giving three arrays and the number of threads to be used as N. In this case N needs to be the length of the arrays.

The `__global__` keyword in front of the function declaration tells nvcc that it is a kernel. The `<<< N, M >>>` specifier tells the compiler how many N blocks and how many M threads per block that the kernel should use[29]. A block is a group of threads and shares some memory and resources. The blocks can be executed in any order so they need to be completely independent from each other but variables can be shared among threads inside a block[29]. This allows the program to scale to different GPUs as shown in figure 2.10. The blocks are organized in a one, two or three dimensional *grid*. These can be accessed by each thread so it knows which grid it is in as illustrated in figure 2.11. This can be used to make it easier to assign the threads to the correct bit of the calculation. For instance using a two dimensional grid is good for matrices[29, 33].

Example 3: Example of a declaration of a simple kernel

Data: Vectors A, B, C as C-style arrays

```
__global__ void Add(float* A, float* B, float* C)
{
    int i = threadIdx.x;
    C[i] = A[i] + B[i];
}
```

Example 4: Host code to call the kernel in example 3 with N threads

Data: Vectors A , B , C of length N as C-style arrays

Add<<<1,N >>>(A, B, C);

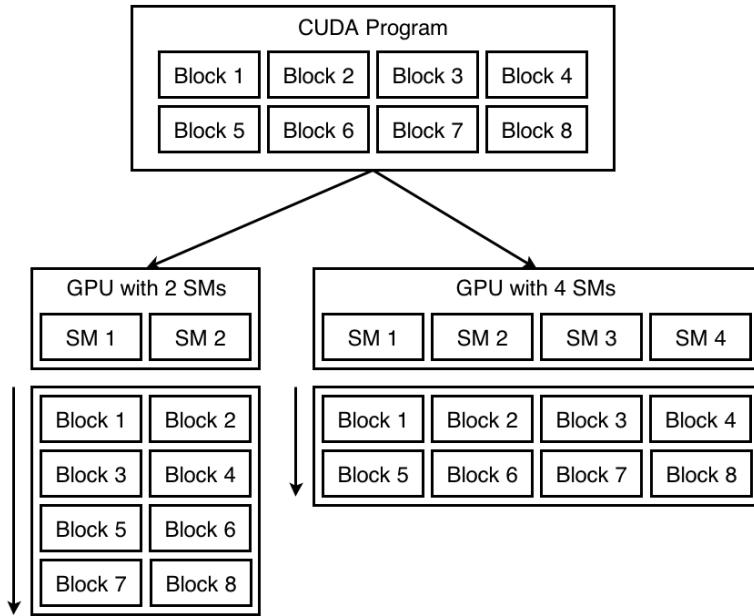


Figure 2.10: Blocks execution depending on the number of streaming multiprocessors

The streaming multiprocessors(SM) is the hardware that handles the execution of these blocks and can execute hundreds of threads concurrently. It uses *SIMT*(Single-Instruction, Multiple-Thread) which is similar to SIMD described earlier in section 2.2.2. SIMT works mostly as SIMD except that it can act as MIMD on collections of threads called *warps*. Each warp consists of 32 threads. When an SM gets a block it is split into warps that are assigned to warp schedulers. Each warp scheduler gives one instruction to a warp so full efficiency is achieved when all the 32 threads perform the same instruction. If there is any divergence it has to disable unrelated threads, so divergence can be costly. However different groups of warps are on different warps schedulers so can diverge without problem.[29]

2.3.1 Device Memory

The GPUs memory is physically on a different device separate from the computers main memory which means that they have separate memory spaces. An object in one memory is not accessible in the other memory. The computers main memory is the *host memory* and the GPUs memory is the *device memory*. Since they are separate data has to be transferred to the device memory and this is done by explicit calls to transfer sections of the hosts memory.[29]

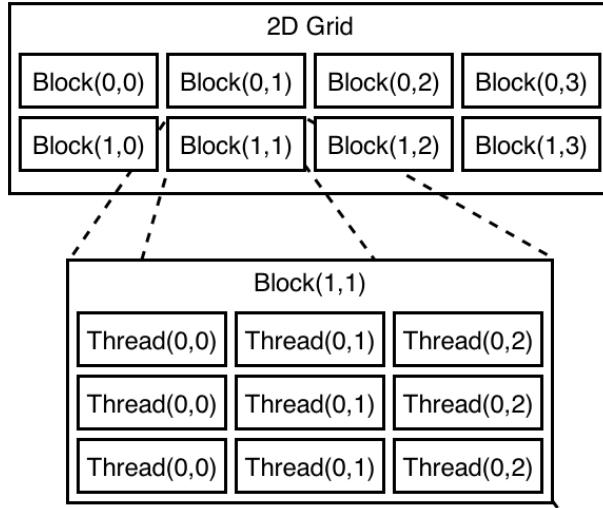


Figure 2.11: 2D grid of blocks

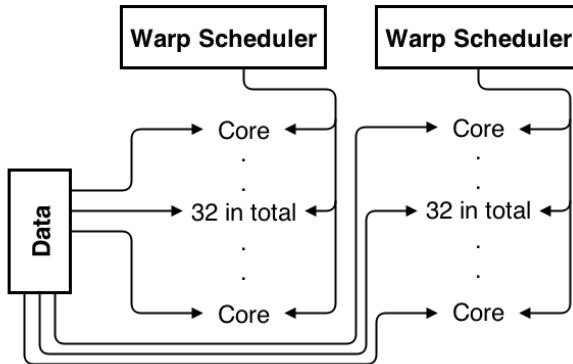


Figure 2.12: SIMT architecture

The GPU also have several different types of memory[29]. Correct usage can give increased speed[29, 33].

- Register memory is located on the multiprocessor and usually costs zero cycles to access. The multiprocessor splits the available registers over its threads so if there are many threads that uses many variables not all of them will fit in the register. This is why a program sometimes can be faster with lower number of threads.[29]
- Global memory is the main memory of the GPU and is accessible from all threads and blocks. However it is relatively slow to access.[29]
- Shared memory is shared inside a block and is faster than global memory. However it is limited in size.[29]
- Constant memory is small however it is read only which enables some optimizations. It is best used for small variables that all threads access.[29]

- Local memory is tied to the threads scope, however it still resides off-chip so it has the same access time as global memory.[29]
- Texture memory is read only and can be faster to access than global memory in some situations. This was more important in older GPUs when global memory was not cached.[40, 29]
- Read-only cache is available on GPUs based on Kepler architecture and uses the same cache as the texture memory. The data as to be read only each multiprocessor can have up to 48kb of space depending on GPU.[43]

2.3.2 Streams

The kernel calls can be made on a *stream*. The easiest way to think of the stream is as queue. The kernels on a stream will be executed in the order they are made but kernels from different streams can be executed in any order given there is enough computational resources. In this way it is possible to execute up to 32 concurrent kernels depending on GPU.[29]

This is not entirely true for some GPU architectures because they only have one queue for the kernels. Because of this independent kernels from different streams can block each other. If we have a group of kernels and issue them on stream one first and then on stream two all the kernels in stream two will be stuck in the queue waiting for the kernels in stream one to finish. This is because the kernels queue only looks at the first kernels to see if they can be executed. The second kernel in the first stream then blocks the queue from moving forward so it will not see the kernels on stream two that could be executed. Newer architectures after Kepler on the other hand has several queues so they don't have this problem.[29, 33, 43]

Streams are also important for asynchronous transfers. These transfers are executed on a stream and just as kernels gets executed after the previous kernels on the same stream is done. The advantage is that other streams can do calculations as normal while the transfer happens. This can hide the time for transfers completely in some situations as shown in 2.14. However the host memory has to be pinned, pinned memory means that the operative system can not page that memory. Paging is that the operative system stores a part of the memory in another area, usually the disk memory, to save space in the memory. Too much pinned memory can slow down the computer.[29, 44]

However there is problem that asynchronous transfers can cause depending on the GPUs architecture[44, 45]. Some older GPUs only have one copy engine while newer have one for each direction, one to the GPU and one from the GPU[44, 45]. This can affect how the calls should be ordered and using the wrong one can make the performance worse than without using asynchronous transfers[44, 45].

There is an example of this in [45] which illustrates the problem. They have four versions of the same code, a sequential transfer version and three asynchronous transfer versions. Two different GPUs were used, one had one copy engine the other had two. In the asynchronous the data is split over four streams coloured differently in the figure 2.14. Version 1 initiates the calls by looping over the

streams one by one and doing the transfer and kernel calls on that stream before moving on to the next[45]. Version 2 makes all the host to device transfer calls for all streams first, then the kernels and then the device to host transfer call[45]. Version 3 is the same as version 2 but with a dummy even after each kernel[45]. The figure 2.14 shows how the transfers and kernels are executed on the GPU.

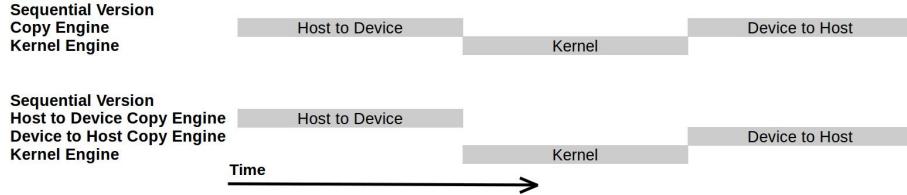


Figure 2.13: Sequential versions

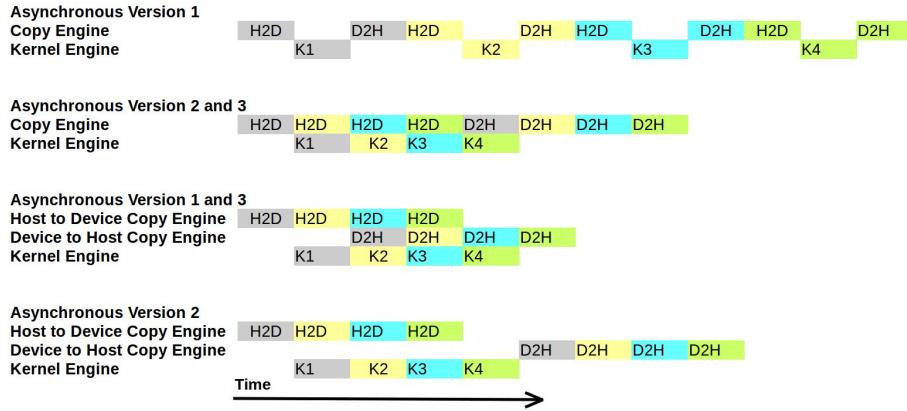


Figure 2.14: Asynchronous versions. D2H=device to host transfer. H2D=host to device transfer. Each colour represents a stream.

2.3.3 Efficient CUDA

One of the main criticism against GPUs for general computing purposes is that it is hard to get good performance because it requires good knowledge about details of the GPU architecture, especially the memory architecture. This section is a summary of things that can be good to consider for CUDA programs from the CUDA programming guide[29], CUDA best practices guide [33] and a similar list as this one, however slightly outdated now since it is from 2010, from a master thesis about GPU in GWAS[40].

Maximize parallelism

Structure the program and the algorithm in such a way that it is as parallel as possible and overlap the serial parts on CPU with calculations on the GPU.[40, 29]

Minimize transfers between host and device

Moving data between host and device is expensive and should be avoided

if possible. It can be better to run serial parts on the GPU rather than moving the data to the host to do the calculation on the CPU. The bandwidth between host and device is one of the large performance bottlenecks. This can be a problem when the data is too large to fit in the relatively small GPU dram.[29, 33]

Find the optimal number of blocks and threads

There are many things affected by the number of blocks and threads so they should be considered carefully. It is a good idea to parametrize them so that they can be changed for future hardware and varied for optimization. NVIDIA has an occupancy calculator which can be helpful in determining the optimal numbers, however high occupancy does not mean high performance.[29, 33]

The number of blocks should be larger than the number of multiprocessors so that all multiprocessors have at least one block to execute. Having two blocks or more per multiprocessor can be good so that there are blocks that aren't waiting for a `__syncthreads()` that can be executed. However this is not always possible due to shared memory usage and similar.[33]

The number of threads per block should be a multiplier of 32 but minimum 64. It's also important to remember that multiple concurrent blocks can reside on the same multiprocessor. Too large number of threads in a block and parts of the multiprocessor might be idle since there aren't a block small enough to use those threads. Between 128 and 256 threads is a good place to start.[33]

Use streams and asynchronous transfers

By using streams it is possible to overlap memory transfers with calculations as mentioned before. This means that the data for the next batch can be transferred while the current batch is calculated and when it is done it can start calculating on the next batch directly after the current one is done. This can hide the time for transfers completely in some situations. Depending on the time the transfers take versus the time the calculations take this can give significant speedup.[29, 44, 43]

Use the correct memory type and caches

Correct use of caches and memory is important for both CPU[26] and GPU. However it is more complicated on GPU since the caches are smaller and there are several types of memory as mentioned before[29, 33].

Avoid divergence

Each thread in a warp executes the same instruction at the same time so if some of threads diverge the rest will be ideal until they are at the same instruction again. This means it is important to use control structures such as `if` statements carefully to prevent threads from idling.[29, 33]

Avoid memory bank conflicts when using shared memory

Shared memory is divided into equally-sized memory modules called banks that can be accessed at the same time for higher bandwidth. Bank conflicts occur when separate threads access the same bank. On some GPUs it is

fine if all threads access the same bank. Bank conflicts are split into as many conflict-free requests as needed.[29, 33]

Use existing libraries

Instead of writing everything from scratch it is usually a good idea to use already existing libraries. Especially when performance is important and most task are non trivial on GPUs so using an already optimized library is a good idea. Some of the most popular libraries for CUDA are:

- CUBLAS: BLAS implementation for CUDA. BLAS and LAPACK is a standard for a library that provides highly optimized functions for linear algebra.[35]
- CULAtools: BLAS and LAPACK implementation for CUDA for both dense and sparse matrices[46]
- MAGMA: BLAS and LAPACK implementation among other things that can distribute the work on both CPU and GPU[47]
- Thrust: Template based library that tries to emulate C++ standard library[48]

Avoid slow instructions

There are some instructions that can be slow and should be avoided if possible, for instance type conversion, integer division and modulo. If a function is called with a floating point number that might be used as a double and require a conversion. By putting an *f* at the end of the number it is told to be single precision float, for instance 0.5f. In some cases it is possible to use bitwise operations instead which is faster.[29, 33]

Restricted pointers can give increased performance

Aliasing can be a problem as mentioned earlier. By using the `__restrict__` keyword on pointers the compiler can be told that no aliasing will occur, however it is up to the programmer to make sure that is the case or there might be unexpected results. Not using aliasing reduces the number of memory accesses the CPU needs to make. However it increases register pressure so it can have an negative effect on performance.[29]

Use fast math functions if precision isn't needed

There are three versions of the math functions. The double precision version is `func()` while the single precision function is `funcf()`. There is third faster but less accurate version `_funcf()`. The option `-use_fast_math` makes the compiler change all the `funcf()` to `_funcf()`.[33]

Instruction level parallelism can increase speed

Just as for CPUs the GPUs can take advantage of instruction level parallelism. By unrolling loops this can give 2x the speed relatively simple[49].

2.4 Software Design Background

How to design software so it can be maintained over time has been a problem for a long time. Object oriented languages such as C and later Java and C++ arose because of problems with maintaining software. Badly organized and written code will cost productivity in the future when bugs and other problems stack up because of earlier mistakes. Fixing those bugs are likely to be time consuming because it can be hard to find where they originated from. All code gathers problems overtime, however a well designed system will degenerate significantly slower than one that was designed carelessly. The code will also be read by ourselves and others later which means readability is important if a future reader is to understand the code and find the parts they are interested in. [50, 51]

This is also important in science where others might wish to use the tools or repeat the experiments. A recent study[52] found that the repeatability in computer science is low. They used the term reproducibility however repeatability is a more fitting term because they didn't try to reproduce the results from scratch, they tried to repeat it using the same code. Of 513 articles they received the code from 231, 108 failed to build, 21 failed to run properly[52]. The remaining 102 compiled and ran, however they did not check if the results of the program was correct[52].

There are ways to design the program and code in such a way that it is easier to read and maintain. Most of the things described here goes under the development technique called agile development and the concept of clean code and is a summary of some concepts from the book Clean Code[50]. SOLID is an acronym for five principles for object oriented programming and design[50]. When used together they are intended to make programs that are easy to maintain and extend[50]. As already mentioned readability is important. Correct names will make the code explain itself without other documentation, the code itself is the documentation[50].

Initial	Principal	Concept
S	Single responsibility principle	A class should only have a single responsibility
O	Open/closed principle	A class should be open for extension but closed for modification
L	Liskov substitution principle	If S is a subtype of T, then objects of type T may be replaced by objects of type S without altering any of the desirable properties of the program(e.g. work performed)
I	Interface segregation principle	Use several smaller and more specific interfaces instead of one large
D	Dependency inversion principle	Depend on abstractions(e.g. interfaces) not details

Table 2.4: The five SOLID principles, from [50]

Dependency injection is one way to implement dependency inversion principle[50]. Injection is passing the dependency to the dependent object. This is used instead of allowing the dependent object to construct or find the dependency.

2.4.0.1 Basic C++ stuff

2.4.1 Unit Tests and Mocks

Unit testing involves testing the program in small units in isolation. Testing in isolation means that the test should only depend on the part of the program that is tested. If a part of the program is not working properly only its related tests should fail, not other tests for code that depends on it but otherwise works properly. This makes it easier to find errors when they do occur since the tests will pinpoint the unit which does not work. *Integration testing* is used to test several units together to find possible errors

[50]

It can be easy to denote test code as less important than the main code but they should be treated as equally important. Tests should also not be an inconvenience to use so they should be easy to run, take reasonable amount of time to complete and not require any outside interpretation whether they failed or not.[50]

Mocking is to replace a real object with a fake object called a *mock* that for the code is indistinguishable from the real object. This allows one to create situations and test with more control and without depending on the real objects code. The second is important for unit tests since it enables one to test units that normally depend on others in isolation. In the first case it is useful in situations such as when one wants to test a class handling of a rare failure. Such failures can be hard and time consuming to induce. It is then easier to use a mock object that behaves like the failure has occurred. However mocking requires that the code for the class doesn't create the object itself since there is no way to replace the object with the mock. This is one of the reasons why dependency injection should be used.[50]

2.4.2 Design Patterns

A *design pattern* in software design is a reusable general solution to a common recurring problem in a given context. It is templates and structures to solve the problem, however it is not code. However they are partially dependent on the programming language since different languages have different features and limitations.[51]

Consumer producer is a concurrency pattern for when there are a number of consumers/workers and producers. They share a common queue for products. The producer generates some product and put it into the queue, while the consumers consume the products. The problem is that the producers should not add to an non empty slot and that each product should only be consumed once.[50, 51]

2.4.3 Version Control Software

Backups for data, code, text and so on is important if something happens. One way to backup text and code projects is to use a version control software and is usually an important tool in developing[53]. For this project Git was used and the source code can be found at <https://github.com/Berjiz/CuEira>. Version control software allows several developers to share a common *repository* which allows them to work on different parts at the same time, it could even be in the same file[53]. However changes at the same place causes *merge conflicts* which usually needs to be solved manually[53]. A version control software keeps track of all the changes made so if a part later turns out to be wrong that part can reverted to an older correct version[53]. Version control software also have many other features[53].

2.5 Performance Measures

An important part in creating fast and efficient programs is to know how fast the program is under certain conditions and which parts of the program are slow[22]. For instance the speed could suddenly drop when too many threads, there might be a bottleneck in the communication, and so on.

There are two ways to measure how long a program takes to execute[22]. Wall clock time is how long real life time the program took. The other is to measure the number of processor cycles spent. A parallel program will have shorter execution time than it is serial version however it will likely have spent more processor cycles due to overhead from communication and initialization of the threads. These two measures are useful for different kinds of comparisons. Wall clock time is better for overall performance while number of cycles is useful for comparing different algorithms[22, 33].

Speed up is a measure of how much faster then program is with a certain number of threads compared to the serial version. It's defined as[22]

$$S(p) = \frac{T(1)}{T(p)}$$

Where $T(1)$ is execution time of serial program and $T(p)$ is execution time of parallel program with p threads. Linear speed up is when $S(p)=p$ [22].

Efficiency reflects how efficient the program is using p threads. Linear speed has efficiency 1. It's defined as[22]

$$E(p) = \frac{S(p)}{p} = \frac{T(1)}{pT(p)}$$

Strong scaling refers to how the program handles a fixed problem size and increased number of processors[22]. An program with strong scaling has linear speed up[22]. Weak scaling refers to the execution time of the program when there is a fixed problem size *per processor* and the number of processors is increased[22, 33].

It can be a good idea to plot these measures while varying p , this can show when a bottleneck occurs. Looking at the measures at node level can also be useful to get an idea of how increased number of nodes and therefore increased communication over the network affects the performance.

2.5.1 Amdahl's Law and Gustafson's Law

Amdahl's Law is used to find the expected speed up of a system when parts of it are made concurrent. Simply it says that as the number of processors increases the parts that aren't parallel will start taking up more and more of the wall clock time and that the speed up for adding more processors will decrease as more and more processors are added and more time is spent relatively on the non parallel parts[22]. It's closely related to strong scaling.[33, 54]

It says that the expected speed up with F fraction of the code parallel and p threads is[22]

$$S(p) = \frac{1}{(1 - F) + \frac{F}{p}(1 - F)}$$

As the number of threads grow towards infinity $S(p)$ converges on $\frac{1}{1-F}$. If we have 90% of a code parallel then even with infinite number of threads we won't get a better speed up than ten.[54]

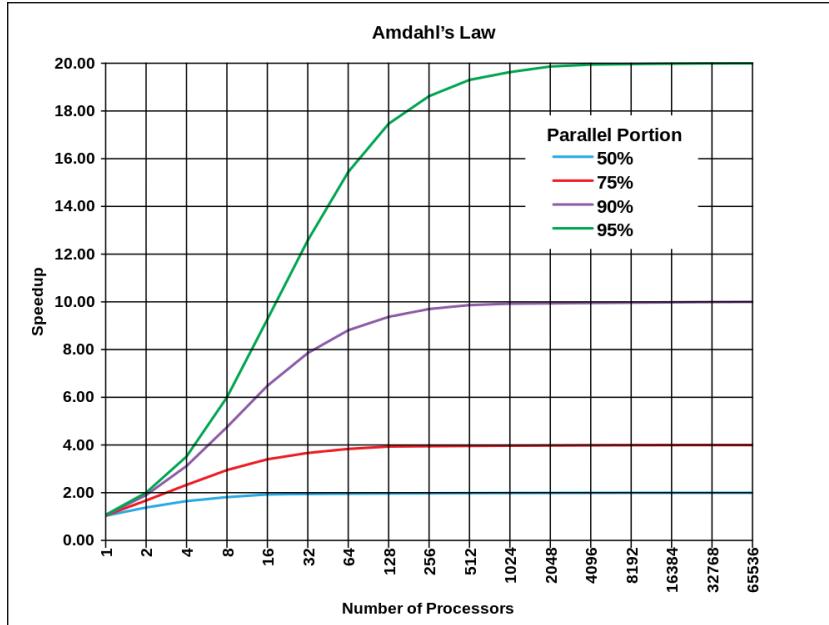


Figure 2.15: Illustration of Amdahl's Law. Wikipedia Commons

There are limitation to Amdahl's Law since it makes a couple of assumptions.

- The number of executing threads remain constant over the course of the program.
- The parallel portion has perfect speed up. Often not true due to shared resources, eg caches, memory bandwidth, and shared data.
- The parallel portion has infinite scaling, not true due to similar limits as above. More threads will not increase performance after a while or might even decrease it.

- There is no overhead for creation and destruction of threads.
- The length of the serial portion is independent of the number of threads. Often the serial work is to divided the work to the threads, this work will obviously increase as the number of threads go up. More threads can also lead to more communication overhead.
- The serial portion can't be overlapped by the parallel parts. For instance with producer consumer type pattern the consumer could be strictly serial but the time it takes could be overlapped by the parallel producers.

This means it is most accurate with programs that are of the fork-join type, e.g. both serial and parallel parts[55].

Gustafson's Law is closely related to Amdahl's Law and can in some ways be more accurate than Amdahl's Law[55]. Gustafson's Law makes similar assumptions as Amdahl's Law however it also makes two additional statements. It states that problems tends to expand when provided with more computational power, e.g. increased precision by reducing grid size for simulations, higher frame rate for graphics and so on[55, 22]. The second is that the parallel portion of the program tends to expand faster than the serial part, e.g. for matrix multiplication the initialization scales linearly with the matrix size while the multiplication itself scales as $O(n^3)$ [55, 22]. The former means that it is closely related to weak scaling[22]. So in a way it says that the execution time remains constant rather than the amount of data. More precise it says that the expected speed up with p threads and F fraction of the code that is parallel is[55, 33, 22]

$$S(p) = p + (1 - F)(1 - p)$$

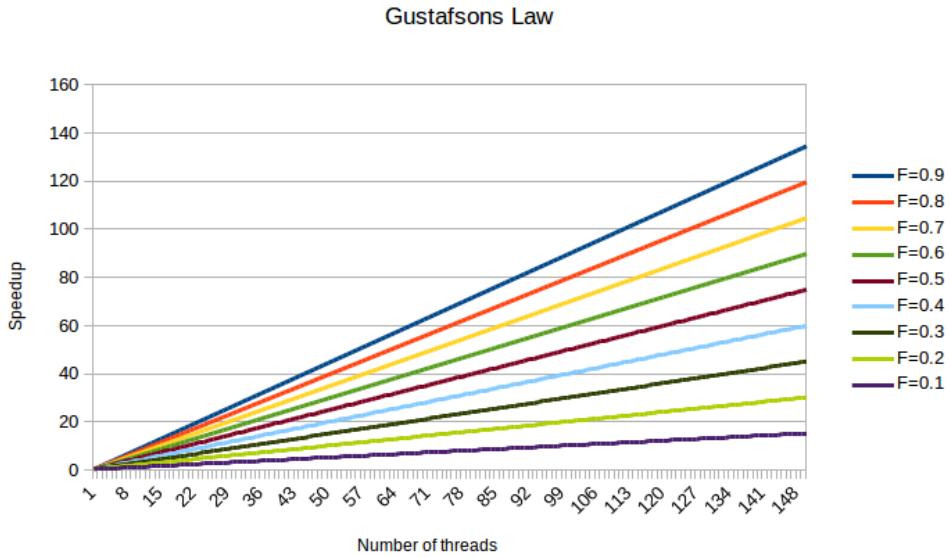


Figure 2.16: Illustration of Gustafson's Law

2.5.2 Profilers

There are applications called profilers that are made to assess the programs performance and resource consumption. They calculate some of the measures mentioned earlier and they also check hardware usage and how much time the program spends at various parts of the program. This is very useful for finding bottlenecks and other problems in the program. It does not matter if the algorithm is super fast if all the data is stuck in network transfers. The profilers can be hardware dependent so the manufacturers usually provided them for their products.[22, 33]

Chapter 3

Algorithms and Implementation

This section starts with the algorithms and then the implementations.

3.1 Algorithmic Background

There are many algorithms and programs proposed for searching for interaction, most have focused on gene-gene interaction as already mentioned[2]. One of the challenges of gene-environment interaction is that environmental factors can be of any variable type(i.e. binary, continuous, categorical) which creates problems in various ways[2]. Gene-gene interaction tools can sometimes be used to find gene-environment interaction, however they usually require the variables to be binary or have other problems since they weren't designed for environmental interaction[2].

Both groups of computers(i.e. clusters) using regular processors[56] and graphic processors[8, 36, 37, 38, 39, 40] have been used for GWAS. Graphic processors for computing have become more popular in the last ten years. They have been a popular choice for some GWAS methods because each combination of variables can commonly be considered independently from the others. More about graphic processors in section 2.2.4 and why the are good for GWAS in section 2.2.4.1.

The methods can be roughly classified into four categories, exhaustive, stochastic, machine learning/data mining and stepwise[10].

Exhaustive search is the most direct approach, it compares all combinations of the SNPs in the dataset. Exhaustive search methods will not miss a significant combination because it didn't consider that specific combination. However it also means that they can be slow since they will spend time on combinations other methods would skip completely. Multifactor-Dimensionality Reduction(MDR)[57] and BOOST[58] are two examples of this type of algorithm.

Stochastic methods uses random sampling to iterate through the data. BEAM[59] is one example and it uses Markov Chain Monte Carlo(MCMC) method.

Data Mining and *Machine Learning* are methods that try to learn patterns from data and tries to generalize it. MDR[57] is a type of data mining method and is among the most common methods used in GWAS. See section 3.1.2 for more details.

Stepwise methods uses a filtering stage and a search stage. At the filtering stage uninteresting combinations are filtered out by using some exhaustive method. The other SNPs are the examined more carefully in the search stage. BOOST[58] is an example which uses succinct data structures and a likelihood ratio test to filter the data before applying log-linear models.

3.1.1 Logistic Regression

One way to model the contingency tables is by using *logistic regression*. Logistic regression is a type of linear regression model for classification that models a latent probability for the outcomes. The outcomes are binary, however the method can be extended to multiple outcomes. In this work we will only consider them as binary. Logistic regression transforms the probability by using the *logit* transformation. The logit transformation with probability π is [17]

$$\log\left(\frac{\pi}{1 - \pi}\right) \quad (3.1)$$

The probability with a set of predictor variables X is $\pi(X) = P(Y = 1)$. The linear regression model with n predictors $X = (x_1, x_2, \dots, x_n)$, coefficients $\beta = (\beta_1, \beta_2, \dots, \beta_n)$ and by using the logit transformation is then[17]

$$\text{logit}[\pi(X)] = \alpha + \beta X \quad (3.2)$$

By moving the logit to the right side of the equation we get the model of the probability[17]

$$\pi(X) = \frac{e^{\alpha+\beta X}}{1 + e^{\alpha+\beta X}} \quad (3.3)$$

The logit, equation 3.1, also happens to be the log of the odds(equation 2.5)[17]. By exponentiating both sides of equation 3.2 it shows that the odds is [17]

$$e^{\text{logit}[\pi(X)]} = \frac{\pi(X)}{1 - \pi(X)} = e^{\alpha+\beta X} = \Omega \quad (3.4)$$

This means that $\exp \beta$ is the odds ratio since the odds increase by $\exp \beta$ for each unit increase of X [17]. It can also been seen by taking the ratio of the odds using equation 2.6 and with $X = x + 1$ and $X = x$

$$\theta = \frac{e^{\alpha+\beta(x+1)}}{e^{\alpha+\beta x}} = e^{\alpha+\beta(x+1)-\alpha-\beta x} = e^\beta \quad (3.5)$$

Finding the β coefficients are done in a similar way as with other linear regression models since they all are generalized linear models [17]. It's usually done using maximum likelihood(ML), via Newtons method[17, 12]. It's an iterative method which means it can be relatively slow compared to non iterative methods. The

pseudo code for the algorithm using Newtons method can be found in algorithm 3 [12]. Line 13 is sometimes inside the loop, however since it is not needed for the actual iteration calculating it for every loop wastes time.

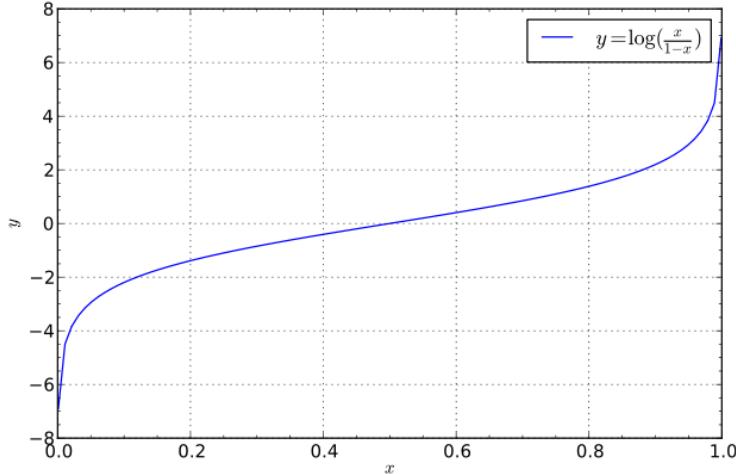


Figure 3.1: Graph of the logit transformation. Wikipedia Commons

3.1.2 Data Mining and Machine Learning Approaches

Approaches based on Data Mining and Machine Learning have been a popular choice for GWAS. MDR[57] and Random Forest(RF)[60] are among the most common ones[2, 1]. There are other methods as well such as clustering approaches [10]. Most of them are used for screening the data for possible interactions[2, 1].

Their biggest advantage is that they are usually non-parametric and designed with high dimensional data in mind. However they are prone to overfitting and the usual way to try to prevent that is to use cross validation and sometimes permutation tests. It means that even if the method itself is fast it is repeated so many times that the whole algorithm can be slow in the end.[1]

3.1.2.1 Multifactor-Dimensionality Reduction

MDR is a method that reduces the number of dimensions(i.e. variable) by combining several dimensions into one. In GWAS it combines the variables that are suspected to interact. This new variable is then compared against the outcome. If the new variables predictability of the outcome is high enough then the variables that were combined are considered to interact. This process is usually repeated on all pair combinations of variables.

The reduction from n dimensions is done by calculating the ratio of cases versus controls for each combination of the possible values of the variables. If the ratio is above a certain threshold all the members of that groups get the value 1 for the new dimension, otherwise 0. Accuracy of the model is done by using cross

Algorithm 3: Logistic regression using maximum likelihood and Newtons method

Data: N number of data points
 M number of variables, excluding the intercept
 \mathbf{X} is an $N \times M$ matrix that contains the variables
 \mathbf{Y} is the outcomes with length N
 β has length M and contains the initial values of β

Note: * is element by element multiplication

```

1  $\mathbf{X} \leftarrow \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$ 
2  $\beta \leftarrow \begin{pmatrix} 0 \\ \beta \end{pmatrix}$ 
3  $iter \leftarrow 0$ 
4  $diff \leftarrow 1$ 
5 while  $iter < max\_iter$  and  $diff > threshold$  do
6    $\beta_{old} \leftarrow \beta$ 
7    $p \leftarrow \frac{e^{\mathbf{X} \cdot \beta}}{1+e^{\mathbf{X} \cdot \beta}}$ 
8    $s \leftarrow \mathbf{X}^T \cdot (\mathbf{Y} - p)$ 
9    $\mathbf{J} \leftarrow (\mathbf{X}^T \cdot (p * (1 - p))) \cdot \mathbf{X}$ 
10   $\beta \leftarrow \beta_{old} + \mathbf{J}^{-1} \cdot s$ 
11   $diff \leftarrow \sum |\beta - \beta_{old}|$ 
12   $iter \leftarrow iter + 1$ 
13  $log\ likelihood \leftarrow \sum(\mathbf{Y} * ln p + (1 - \mathbf{Y}) * ln(1 - p))$ 

```

validation and permutation tests, in simpler words it means that it reshuffles the data randomly and recalculates the model many times to get an estimate of the models certainty. Because of that MDR can be slow. However it is still usually faster than exhaustive search with regression methods.[1, 57]

MDR can been used for gene-environment interaction but requires modifications since MDR can only handle binary variables. There are extensions that can use continues variables, however these are regression based so these will be slower than regular MDR.[2]

A simple example of MDR using exclusive or (XOR). XOR is an logical operator that is true if one and only one of its two variables is true. We have 4 possible combinations and an occurrence for each of them, see table 3.1. The combination (1,0) and (0,1) both have one case with outcome 1 so MDR will classify them as 1 in the new variable Z, the other two combinations have outcome 0 so will be classified with Z=0, see table 3.2. From here it is easy to make an predictor from Z to the outcome Y by comparing the values.

Y	X₁	X₂
1	1	0
1	0	1
0	0	0
0	1	1

Table 3.1: XOR table with outcome **Y** and variables **X₁** and **X₂**.

Y	Z
1	1
1	1
0	0
0	0

Table 3.2: XOR table with **X₁** and **X₂** combined into **Z** using MDR.

3.1.2.2 Random Forest

RF is an ensemble learning method[60]. Ensemble methods combine multiple models to improve performance. RF takes randomized samples of the data and builds decision trees on each of them. These trees are then combined to form the classifier. Usually hundreds or thousands of trees are used depending on the problem[60]. One of the most popular variants of Random Forest for GWAS is Random Jungle[61].

It has been shown in high dimensional data that RF tends to only rank interacting factors high if they have strong marginal effects[62]. Also the ranking of the variables does not indicate which factor it is interacting with either since it is based on the joint distributions[2]. How to incorporate the environmental factors in RF is also not obvious and using variables with very different scales can bias the results[2].

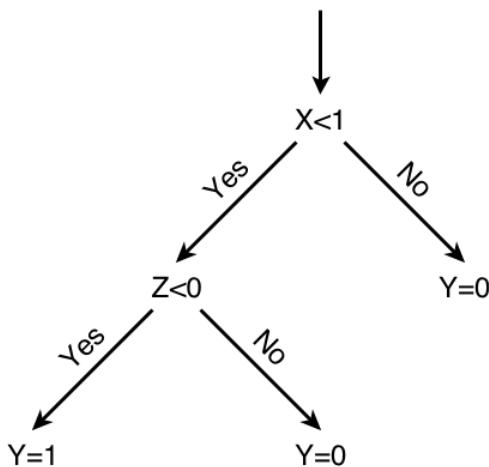


Figure 3.2: A example of a decision tree

3.2 Implementation

Here bla bla
why LR

All the programs read Plink data files except CuEira which can only read the binary Plink files. They also read a separate file that contains the environmental data and covariates if any. For CuEira the covariates and environmental data is split in two different files with the same structure.

3.2.1 JEIRA and GEISA

This section will briefly explain the structure of the GEISA program. Almost everything applies to JEIRA as well. It is written in Java and uses Javas built-in functions for concurrency. Since GEISA is based on JEIRA, as talked about in section 1.5, the underlying design is the same.

Why Java

GEISA has three dependencies, two are provided as JAR files in the distribution so the user doesn't need to get them. Maven is only needed if the user wants to compile the program from the source.

Apache Commons CLI

Provides command line interface, e.g. it parses the options when using the program.

Apache Commons Math

Provides matrix and vector classes, linear algebra, etc.

Apache Maven

Used to compile Java programs.

The implementation uses the producer consumer pattern described in section 2.4.2. The main thread creates a queue of tasks. All the result producers iterates over this shared queue and outputs results. These results are placed in a queue a consumed by another thread that prints the results. The program reads all data at the start which means it can get memory problems for large datasets. GEISA have a few unit tests for some of the basic storage classes and the binary file reader.

3.2.2 CuEira

This section will explain the structure of CuEira and how it works starting with a general overview and how it is compiled.

C++ was chosen because

Why C++ good memory handling with pointers fast blas, cuda commonly C and Fortran interface

CuEira have some dependencies. The dependency on MKL and Intels compiler is something that can be changed. BLAS shares the interface so another BLAS library could work with updates to the CMake file.

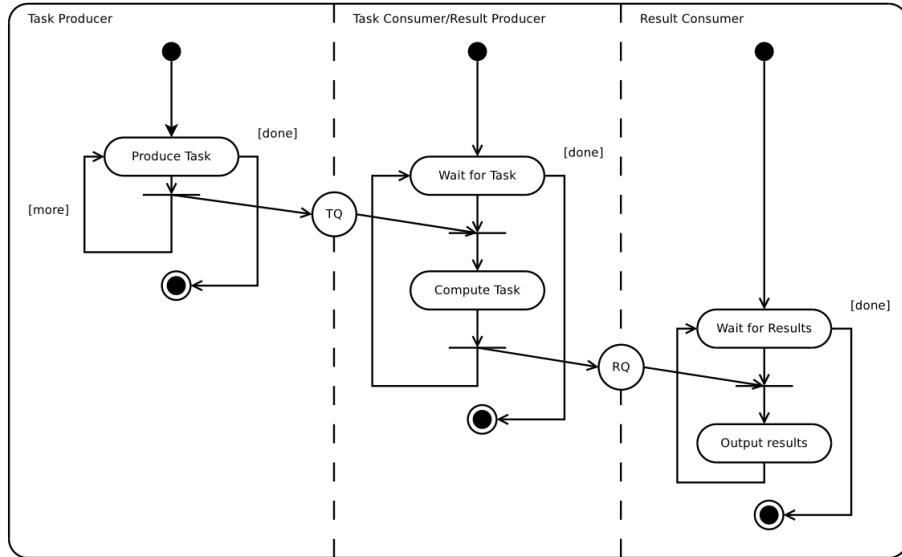


Figure 3.3: JEIRA structure

Make

Used for compiling

CMake

Used to create Make files

Boost

Boost is a collections of libraries

CUDA

NVIDIAs library for GPU, explained in section 2.3

CUBLAS

BLAS for CUDA

MKL

Intels BLAS library

Intel Compiler

Intels C/C++ compiler

Google Test and Google Mock

Unit test and mocking framework. They are provided in the CuEira source so the user does not need to install them separably.

CuEira is compiled by using CMake and Make. CMake generates Make files from a CMakeLists file. Make is

One of the main advantages of CMake is find package

The main CMake file needs to know what and where all the source files are. This is done by using nested CMake files, each directory has an CMakeLists file that adds its sub directories and sets all the source files to the previous level.

There are four types of lists of source files each is compiled into a library and then linked together. There is one list for all the tests, one for the files that needs to be compiled with nvcc, one for other GPU related code and then one with the host code.

By changing an option in the main CMake file the program can be compiled with either single and double precision. More specific details on how to compile CuEira can be found in appendix D.

The program itself is split into two stages, initialisation and calculations. The initialisation part reads all the information except the genetic data of the individuals and creates some of the classes. In the calculation step the worker threads are started and takes a SNP from the queue, reads its data and calculates the model.

Avoiding copy

Most classes have copy and move operator deleted to prevent mistakes since the default copy is naive and therefore usually flawed when dealing with classes with pointers.

3.2.2.1 Namespaces

Namespaces are used in CuEira to

3.2.2.2 Storage

A set of container classes holds the data. There are two types of them, the basic matrix and vector classes and data containers. The data containers are responsible for handling the data for the model and performing the recoding. They use the vector classes for the actual storage.

There are three types of matrix and vector classes which are all stored in column major format. The set with the Device prefix is for storage on the device and the other two are for storage on the host and shares a common interface. The difference between the host sets is that the ones with Pinned prefix uses pinned memory while Regular does not.

Other basic information such as the information about individuals, environmental factors and SNPs are stored in simple classes where each individual, environmental factor or SNP is a single object. *Enums* are used to store some of them. An enum is a data type consisting of a set of named values. This is useful when the variable has a distinct set of values that makes more sense as text rather than something else. For instance in CuEira the individuals gender is stored as enum with the two values MALE and FEMALE. The named enum values are commonly written in only capital letters. The gender could be stored as an int, however MALE is easier to read than an int value and it also avoids potential errors with incorrect values. An enum can only have one of its possible values while an int is not restricted in the same way.

data vectors

DataHandler

EnvironmentFactorHandler

The SNPs are stored in a vector in the class DataQueue. DataQueue provides a single function that removes and returns a SNP from the vector. It has a lock so it can be shared among all worker threads.

3.2.2.3 File I/O

The

CSVReader can read csv files, both the file with the covariates and the file with the environmental data are csv files.

The EnvironmentCSVReader inherits CSVReader, its purpose is to convert the data matrix to EnvironmentFactorHandler.

There is just one writer, the ResultWriter, as its name suggest its responsibility is to write the results. The results are written immediately and then discarded to save memory. There is an lock on the writing so several threads can share it. The results will therefore be written in the order they are completed.

3.2.2.4 Recoding

DataHandler has the responsibility of recoding and delegates it to the data containers.

3.2.2.5 Model

intro

easy to extend

4 types

Model ModelConfiguration ModelResult ModelHandler

Also ModelInformation CombinedResults

Example of extending as basic ExampleModel

Pictures

3.2.2.6 Wrappers

CUDA and BLAS functions and interfaces is in C style. For this reason they were wrapped in classes that are C++ style. These wrappers do the actual function calls to CUDA and BLAS. Since they contain a lot of different functions this was done when a function was needed to save time.

MKL, e.g. Intels BLAS is wrapped in MKLWrapper. The BLAS interface standard is not object oriented which makes

Matrix vector multiplication with MKL is shown in 5. The function call uses 12 parameters, most of which are of similar types. As is sometimes common it also has gone a bit too far in the shortening of the function name. The *s* stands for single precision, *ge* for general, *m* for matrix and *v* for vector. General means that it is standard matrices and vectors, e.g. not triangular or banded etc. The equivalent function in MKLWrapper can be found in 6. The name matrixVectorMultiply instantly tells the user what it does. As mentioned earlier in section 3.2.2 the CuEira is either single or double precision depending on an option when compiling so MKLWrapper doesn't mention the precision. All matrices and vectors are also general. However even if the function name would contain that information too it would still not be so long that it was in the way.

A possible name and by using namespaces could be:
General::SinglePrecision::matrixVectorMultiply

One time during the thesis it took several days to find an error that was in one of the CUBLAS calls. Due to a 1 instead of 0 it performed the operation $y = z * x + y$ instead of $y = z * x$. These types of errors are much easier to make when a function has a larger number of similar parameters.

Example 5: Single precision matrix vector multiplication using MKL,
 $result = \alpha * matrix * vector + \beta * result$

Data: *layout*, storage type of the matrix, either column major(CblasColMajor) or row major(CblasRowMajor)
trans, use transpose, conjugate or neither on the matrix
matrix, *vector* and *result* are float pointers
 α and β are constants
m, number of rows in the matrix
n, number of columns in the matrix
ld_matrix, leading dimension of the matrix
inc_vector and *inc_result* are the increments of the elements

```
cblas_sgemv(layout, trans, m, n, alpha, matrix, ld_matrix, vector,
inc_vector, beta, result, inc_result);
```

Example 6: Matrix vector multiplication using MKLWrapper, $result = \alpha * matrix * vector + \beta * result$

Data: *matrix*, HostMatrix
vector, HostVector
result, HostVector
 α , constant, default 1
 β , constant, default 0

```
matrixVectorMultiply(matrix, vector, result, alpha, beta);
```

CUBLAS is wrapped together with the created kernels in KernelWrapper in a similar way to how MKL is wrapped in MKLWrapper. KernelWrapper is explained in 3.2.3 Other CUBLAS utility and CUDA functions are wrapped in CudaAdapter.

Stream class

Device class

The transfers to and from the device is handle by two classes, DeviceToHost and HostToDevice. They perform asynchronous transfers on the corresponding matrix and vector class. It can either create a new container to transfer to or transfer to an given place in the memory. The former is useful for instance when transferring different vectors to into a combined matrix.

3.2.3 Kernels

kernelwrapper

same structure, all do operation on elements of vector

Algorithm 4: Kernel for vector addition

Data: Asdf

```
__global__ void ElementWiseAddition(const PRECISION* vector1,
const PRECISION* vector2, PRECISION* result, const int length)
{
    int threadId = blockDim.x * blockIdx.x + threadIdx.x;

    if (threadId < length)
    {
        result[threadId] = vector1[threadId] + vector2[threadId];
    }
}
```

3.2.4 GPU stuff

3 threads per GPU Each thread has one CUDA stream Each thread gets a snp from the queue, reads in the snp data, calculates the models and writes the results

Outcomes shared over whole GPU

3.2.4.1 Logistic Regression

Structure things, inheritance

Some parts on GPU, some parts on CPU Lines 10 and 11 is done on CPU in algorithm 3

The CPU parts of GPU LR is shared by using inheritance with CPU LR

This is because two reasons, the size is small and CUBLAS does not provided an SVD implementation. CULA tools and MAGMA does but CULA tools synchronises over the whole device and MAGMA uses both CPU and GPU is optimized for large matrices.

Other parts are done by using CUBLAS or custom kernels with elementwsie operations.

syncing

3.2.4.2 Memory Usage

The memory usage is low because only the SNPs whose models currently is being calculated is stored in memory. This also means that most of the memory usage depends on the number of individuals not the number of SNPs especially since the SNP class is small. The memory is also bound by the GPU rather than host memory because of the relative size and that the GPU and host need about the same about of memory. However it does depend on the number of GPUs. A Tesla GPU has 6GB of memory[32], how many individuals can it fit

at most?

Float=4 bytes, double=8 With n number of individuals and m , $m = 4 + \text{number of covariates}$, one set of memory for the LR model needs:

The outcomes, size n , shared for all streams

Predictors, $n*m$ Probabilities n Score m Beta m Work area n Work area $n*m$
Information matrix $m*m$

With three streams the number of floats needed is

$$6GB = (n + 3(2 * n * m + 2 * n + 2 * m + m^2)) * precision \quad (3.6)$$

Breaking out n

$$n = \frac{\frac{6GB}{6m+7} - m^2 - 6m}{precision} \quad (3.7)$$

See table 3.3. Even with 100 covariates and double precision it can store 1.2 million individuals.

Number of covariates	Single precision	Double precision
0	48.4	24.2
5	24.6	12.3
10	16.5	8.2
20	9.0	5.0
100	2.4	1.2

Table 3.3: The number of individuals in millions needed to fill 6GB with the given precision, number of covariates and three streams.

n individuals s snps e env c cov

shared outcomes s snps e environment factors $n*e$ storage of env $n*c$ storage of cov

per thread $2*n$ snpvector n envfactor n interaction

3.2.4.3 Testing

Almost all classes have a corresponding unit test that tests its functionality in isolation. However the matrix and vector classes are never mocked out because they are too basic and fundamental. It does mean that if there are multiple errors when testing and the tests for the matrix and vectors failed the fault likely is with those. However the number of integration tests are few and there is a need to create more.

Since mocks are created by inheritance from the class to be mocked all functions need to be virtual. Some classes also have empty protected constructors to be used by the mock, others have helper functions to construct the mock. This goes against some parts of clean code but is necessary because of how C++ and mocking works. It could potentially be improved by using private constructors and friends to give the mock access.

Chapter 4

Results

Intro

To make the comparison of the performance of the programs more fair a couple of details in CuEira was changed. With the changes CuEira assigns the risk allele the same way as the older programs and doesn't use recoding on the multiplicative model. The MAF threshold was also set to zero. The way the risk allele is calculated in GEIRA, JEIRA and GEISA is not according to the definition earlier in section 1.3.1, the actual way it is done is shown in appendix C.1.

Chapter 5

Discussion and Conclusions

Chapter 6

Outlook

Bibliography

- [1] H. J. Cordell, “Detecting gene–gene interactions that underlie human diseases,” *Nature Reviews Genetics*, vol. 10, no. 6, pp. 392–404, 2009.
- [2] S. J. Winham and J. M. Biernacka, “Gene–environment interactions in genome-wide association studies: current approaches and new directions,” *Journal of Child Psychology and Psychiatry*, vol. 54, no. 10, pp. 1120–1134, 2013.
- [3] B. Ding, H. Källberg, L. Klareskog, L. Padyukov, and L. Alfredsson, “Geira: gene-environment and gene-gene interaction research application,” *European Journal of Epidemiology*, vol. 26, no. 7, pp. 557–561, 2011.
- [4] H. Mahdi, B. A. Fisher, H. Källberg, D. Plant, V. Malmström, J. Rönnelid, P. Charles, B. Ding, L. Alfredsson, L. Padyukov, *et al.*, “Specific interaction between genotype, smoking and autoimmunity to citrullinated α -enolase in the etiology of rheumatoid arthritis,” *Nature genetics*, vol. 41, no. 12, pp. 1319–1324, 2009.
- [5] K. J. Rothman, S. Greenland, and T. L. Lash, *Modern epidemiology*. Lippincott Williams & Wilkins, 2008.
- [6] C. Mann, “Observational research methods. research design ii: cohort, cross sectional, and case-control studies,” *Emergency Medicine Journal*, vol. 20, no. 1, pp. 54–60, 2003.
- [7] P. R. Burton, D. G. Clayton, L. R. Cardon, N. Craddock, P. Deloukas, A. Duncanson, D. P. Kwiatkowski, M. I. McCarthy, W. H. Ouwehand, N. J. Samani, *et al.*, “Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls,” *Nature*, vol. 447, no. 7145, pp. 661–678, 2007.
- [8] B. Goudey, D. Rawlinson, Q. Wang, F. Shi, H. Ferra, R. M. Campbell, L. Stern, M. T. Inouye, C. S. Ong, and A. Kowalczyk, “Gwis-model-free, fast and exhaustive search for epistatic interactions in case-control gwas,” *BMC genomics*, vol. 14, no. Suppl 3, p. S10, 2013.
- [9] G. Fang, M. Haznadar, W. Wang, H. Yu, M. Steinbach, T. R. Church, W. S. Oetting, B. Van Ness, and V. Kumar, “High-order snp combinations associated with complex diseases: efficient discovery, statistical power and functional interactions,” *PloS one*, vol. 7, no. 4, p. e33531, 2012.

- [10] S. Leem, H.-h. Jeong, J. Lee, K. Wee, and K.-A. Sohn, “Fast detection of high-order epistatic interactions in genome-wide association studies using information theoretic measure,” *Computational Biology and Chemistry*, 2014.
- [11] D. Sadava, D. Hillis, C. Heller, G. Orians, and W. Purves, *Life The Science of Biology*. Sinauer Associates, 8th ed., 2008.
- [12] D. Uvehag, “Design and implementation of a computational platform and a parallelized interaction analysis for large scale genomics data in multiple sclerosis,” 2013.
- [13] D. Uvhage and H. Zazzi, “Geisa.” <https://github.com/menzzana/geisa>, 2014.
- [14] A. Ahlbom and L. Alfredsson, “Interaction: a word with two meanings creates confusion,” *European journal of epidemiology*, vol. 20, no. 7, pp. 563–564, 2005.
- [15] K. J. Rothman, *Epidemiology: an introduction*. Oxford University Press, 2002.
- [16] A. Sjölander, W. Lee, H. Källberg, and Y. Pawitan, “Bounds on causal interactions for binary outcomes,” *Biometrics*, 2014.
- [17] A. Agresti, *Categorical data analysis*. John Wiley & Sons, Second ed., 2002.
- [18] B. Lindgren, *Statistical theory*. CRC Press, fourth ed., 1993.
- [19] H. T. O. Davies, I. K. Crombie, and M. Tavakoli, “When can odds ratios mislead?,” *Bmj*, vol. 316, no. 7136, pp. 989–991, 1998.
- [20] M. J. Knol, T. J. VanderWeele, R. H. Groenwold, O. H. Klungel, M. M. Rovers, and D. E. Grobbee, “Estimating measures of interaction on an additive scale for preventive exposures,” *European journal of epidemiology*, vol. 26, no. 6, pp. 433–438, 2011.
- [21] J. M. Bland and D. G. Altman, “Multiple significance tests: the bonferroni method,” *Bmj*, vol. 310, no. 6973, p. 170, 1995.
- [22] G. Hager and G. Wellein, *Introduction to high performance computing for scientists and engineers*. CRC Press, 2010.
- [23] Intel, “Intel unleashes its first 8-core desktop processor.” http://newsroom.intel.com/community/intel_newsroom/blog/2014/08/29/intel-unleashes-its-first-8-core-desktop-processor, 2014.
- [24] C. Hoare, *Communicating sequential processes*. Prentice-Hall, Inc., 1985.
- [25] W. F. Gilreath and P. A. Laplante, *Computer Architecture: A Minimalist Perspective: Dynamics and Sustainability*. Springer, 2003.
- [26] U. Drepper, “What every programmer should know about memory,” 2007.

- [27] F. Abi-Chahla, “Intel core i7 (nehalem): Architecture by amd?” <http://www.tomshardware.com/reviews/Intel-i7-nehalem-cpu,2041-2.html>, 2008.
- [28] Intel, “i7-5960x specification sheet.” <http://ark.intel.com/products/82930>, 2014.
- [29] NVIDIA, *NVIDIA CUDA C Programming Guide*, v5.5 ed., July 2013.
- [30] J. Jeffers and J. Reinders, *Intel Xeon Phi coprocessor high performance programming*. Newnes, 2013.
- [31] M. H. NVIDIA, “Maxwell: The Most Advanced CUDA GPU Ever Made.” <http://devblogs.nvidia.com/parallelforall/maxwell-most-advanced-cuda-gpu-ever-made/>, 2014.
- [32] NVIDIA, “Nvidia tesla-kepler product description.” <http://www.nvidia.com/content/tesla/pdf/NVIDIA-Tesla-Kepler-Family-Datasheet.pdf>.
- [33] NVIDIA, *NVIDIA CUDA C Best Practices Guide*, v5.5 ed., July 2013.
- [34] Intel, “Intel Xeon Phi Coprocessor DataSheet, Document ID 328209 003EN.” <http://www.intel.com/content/www/us/en/processors/xeon/xeon-phi-coprocessor-datasheet.html>, 2014.
- [35] NVIDIA, *CUBLAS Library User Guide*, v5.5 ed., July 2013.
- [36] L. S. Yung, C. Yang, X. Wan, and W. Yu, “Gboost: a gpu-based tool for detecting gene–gene interactions in genome–wide case control studies,” *Bioinformatics*, vol. 27, no. 9, pp. 1309–1310, 2011.
- [37] Z. Zhu, X. Tong, Z. Zhu, M. Liang, W. Cui, K. Su, M. D. Li, and J. Zhu, “Development of gmdr-gpu for gene-gene interaction analysis and its application to wtccc gwas data for type 2 diabetes,” *PloS one*, vol. 8, no. 4, p. e61943, 2013.
- [38] S. Lee, M.-S. Kwon, I.-S. Huh, and T. Park, “Cuda-lr: Cuda-accelerated logistic regression analysis tool for gene-gene interaction for genome-wide association study,” in *Bioinformatics and Biomedicine Workshops (BIBMW), 2011 IEEE International Conference on*, pp. 691–695, IEEE, 2011.
- [39] S. Chikkagoudar, K. Wang, and M. Li, “Genie: a software package for gene-gene interaction analysis in genetic association studies using multiple gpu or cpu cores,” *BMC research notes*, vol. 4, no. 1, p. 158, 2011.
- [40] J. Poznanovic, “Plinkgpu: A framework for gpu acceleration of whole genome data analysis,” Master’s thesis, School of Informatics, University of Edinburgh, 2010.
- [41] R. Jiang, F. Zeng, W. Zhang, X. Wu, and Z. Yu, “Accelerating genome-wide association studies using cuda compatible graphics processing units,” pp. 70–76, 2009.

- [42] S. Almeida, “An introduction to high performance computing,” *International Journal of Modern Physics A*, vol. 28, no. 22n23, p. 1340021, 2013.
- [43] NVIDIA, “Kepler Tuning Guide.” <http://docs.nvidia.com/cuda/kepler-tuning-guide>.
- [44] Mark Harris, “How to Overlap Data Transfers in CUDA C/C++.” <http://devblogs.nvidia.com/parallelforall/how-overlap-data-transfers-cuda-cc/>.
- [45] G. Ruetsch and B. Leback, “CUDA Fortran Asynchronous Data Transfers.” <http://www.pgroup.com/lit/articles/insider/v3n1a4.htm>.
- [46] “CULA Tools: GPU Accelerated Linear Algebra,” 2010.
- [47] S. Tomov, R. Nath, H. Ltaief, and J. Dongarra, “Dense linear algebra solvers for multicore with gpu accelerators,” in *Parallel & Distributed Processing, Workshops and Phd Forum (IPDPSW), 2010 IEEE International Symposium on*, pp. 1–8, IEEE, 2010.
- [48] J. Hoberock and N. Bell, “Thrust: A parallel template library,” 2010. Version 1.7.0.
- [49] V. Volkov, “Unrolling parallel loops,” *Tutorial at the*, p. 133, 2011.
- [50] R. C. Martin, *Clean code: a handbook of agile software craftsmanship*. Pearson Education, 2008.
- [51] R. C. Martin, “Design principles and design patterns.” http://www.objectmentor.com/resources/articles/Principles_and_Patterns.pdf, 2000.
- [52] G. Moraila, A. Shankaran, Z. Shi, and A. M. Warren, “Measuring reproducibility in computer systems research.” <http://reproducibility.cs.arizona.edu/v1/tr.pdf>, 2013.
- [53] J. Loeliger and M. McCullough, *Version Control with Git: Powerful tools and techniques for collaborative software development*. "O'Reilly Media, Inc.", 2012.
- [54] X.-H. Sun and Y. Chen, “Reevaluating amdahlâŽs law in the multicore era,” *Journal of Parallel and Distributed Computing*, vol. 70, no. 2, pp. 183–188, 2010.
- [55] J. L. Gustafson, “Reevaluating amdahl’s law,” *Communications of the ACM*, vol. 31, no. 5, pp. 532–533, 1988.
- [56] A. Gyenesei, J. Moody, A. Laiho, C. A. Semple, C. S. Haley, and W.-H. Wei, “Biforce toolbox: powerful high-throughput computational analysis of gene–gene interactions in genome-wide association studies,” *Nucleic acids research*, vol. 40, no. W1, pp. W628–W632, 2012.

- [57] M. D. Ritchie, L. W. Hahn, N. Roodi, L. R. Bailey, W. D. Dupont, F. F. Parl, and J. H. Moore, “Multifactor-dimensionality reduction reveals high-order interactions among estrogen-metabolism genes in sporadic breast cancer,” *The American Journal of Human Genetics*, vol. 69, no. 1, pp. 138–147, 2001.
- [58] X. Wan, C. Yang, Q. Yang, H. Xue, X. Fan, N. L. Tang, and W. Yu, “Boost: A fast approach to detecting gene-gene interactions in genome-wide case-control studies,” *The American Journal of Human Genetics*, vol. 87, no. 3, pp. 325–340, 2010.
- [59] Y. Zhang and J. S. Liu, “Bayesian inference of epistatic interactions in case-control studies,” *Nature genetics*, vol. 39, no. 9, pp. 1167–1173, 2007.
- [60] L. Breiman, “Random forests,” *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [61] D. F. Schwarz, I. R. König, and A. Ziegler, “On safari to random jungle: a fast implementation of random forests for high-dimensional data,” *Bioinformatics*, vol. 26, no. 14, pp. 1752–1758, 2010.
- [62] S. J. Winham, C. L. Colby, R. R. Freimuth, X. Wang, M. de Andrade, M. Huebner, and J. M. Biernacka, “Snp interaction detection with random forests in high-dimensional genetic data,” *BMC bioinformatics*, vol. 13, no. 1, p. 164, 2012.

Appendices

Appendix A

Lists

List of Variables

θ_{ij} Odds ratio with exposure levels i and j

RR_{ij} Relative risk with exposure levels i and j

X Predictors

Y The outcome, binary variable

$\pi(x) = P(Y = 1)$

Ω Odds

List of Abbreviations

MAF

Minor allele frequency, the frequency of the least common allele.

RF

Random Forest

MDR

Random Forest

OR

Random Forest

RR

Random Forest

RERI

Random Forest

AP

Random Forest

SI

Random Forest

GPU

Random Forest

CPU

Random Forest

SM

Streaming multiprocessor

SIMD

Random Forest

MIMD

Random Forest

SIMT

Random Forest

SISD

Random Forest

LR

Random Forest

MIDSD

Random Forest

List of Algorithms

2.1 Basic algorithm for matrix matrix multiplication of two $N \times N$ matrices	13
2.2 Matrix matrix multiplication optimized	13
3.3 Logistic regression using maximum likelihood and Newtons method . .	37
C.4GEIRA pseudo code to determine the risk allele	62
C.5JEIRA and GEISA pseudo code to determine the risk allele	63

List of Examples

2.1 Independent Instructions	14
2.2 Pointer Aliasing	15
2.3 Example of a declaration of a simple kernel	20
2.4 Host code to call the kernel in example 3 with N threads	21
3.5 Single precision matrix vector multiplication using MKL, $result = \alpha * matrix * vector + \beta * result$	43
3.6 Matrix vector multiplication using MKLWrapper, $result = \alpha * matrix * vector + \beta * result$	43

List of Figures

2.1	Illustration of a simple case of confounding. If we do not observe Z we might falsely find an association between X and Y	8
2.2	The layout of the Haswell architecture i7-5960x, from [23]	9
2.3	Illustration of the dining philosophers problem. Wikipedia Commons	10
2.4	Flynns taxonomy of parallel architectures[22]	11
2.5	Layout of a CPU core, Intel Nehalem i7, from [27]	11
2.6	Theoretical throughput of some NVIDIA GPUs and Intel Processors, from[29]. However the chart does not contain a Xeon Phi processor.	15
2.7	Layout of the Maxwell architecture, from [31]	18
2.8	Knights Corner silicon layout, from[34]	19
2.9	Basic overview of a cluster	19
2.10	Blocks execution depending on the number of streaming multiprocessors	21
2.11	2D grid of blocks	22
2.12	SIMT architecture	22
2.13	Sequential versions	24
2.14	Asynchronous versions. D2H=device to host transfer. H2D=host to device transfer. Each colour represents a stream.	24
2.15	Illustration of Amdahl's Law. Wikipedia Commons	31
2.16	Illustration of Gustafson's Law	32
3.1	Graph of the logit transformation. Wikipedia Commons	36
3.2	A example of a decision tree	38
3.3	JEIRA structure	40

List of Tables

1.1	The phenotype based on the genetic model and if the allele with the effect is A	2
2.1	Contingency table describing the outcome of a study, from [17], page 42	5
2.2	Coding of the variables for LR	7
2.3	Cases of recoding	7
2.4	The five SOLID principles, from [50]	27
3.1	XOR table with outcome Y and variables X₁ and X₂	38
3.2	XOR table with X₁ and X₂ combined into Z using MDR.	38
3.3	The number of individuals in millions needed to fill 6GB with the given precision, number of covariates and three streams.	45

Appendix B

File Formats

B.1 Plink Binary Data Format

B.2 Environmental and Covariates File Format

The environmental factors and covariates are stored in separate files with the same format, one file for environmental factors and one for the covariates if any. One of the columns needs to contain the individual ids, the rest of the columns should be data. The delimiter can be anything of the usual ones and there is an option to provided it to the program. The default is tab delimited.

Appendix C

Bugs and Issues

C.1 Risk Allele Code from GEIRA, JEIRA and GEISA

The code for how GEIRA, JEIRA and GEISA calculates the risk allele is not the same as it is define in [3, 12]. The risk allele is different when the most common allele in case and control is different

Algorithm 5: GEIRA pseudo code to determine the risk allele

Data: caseMax_n and controlMax_n is the frequency for the cases and controls.
caseMax and controlMax is the name of the allele that has the highest frequency in cases and controls.

```
if  $\text{caseMax}_n > \text{controlMax}_n$  and  $\text{caseMax} = \text{controlMax}$  then
    riskAllele  $\leftarrow$  caseMax
else if  $\text{caseMax}_n \leq \text{controlMax}_n$  and  $\text{caseMin} = \text{controlMin}$  then
    riskAllele  $\leftarrow$  caseMin
else if  $\text{caseMax}_n > \text{controlMax}_n$  and  $\text{caseMax} \neq \text{controlMax}$  then
    riskAllele  $\leftarrow$  caseMax
else if  $\text{caseMax}_n \leq \text{controlMax}_n$  and  $\text{caseMin} \neq \text{controlMin}$  then
    riskAllele  $\leftarrow$  caseMin
```

C.2 Encountered Bugs

This is a small collection of various bugs and limitations encountered in the course of the project that can be good to know for anyone that wants to compile CuEira or if the libraries are used in other projects.

There are several problems with Intel compilers and libraries. Most libraries work fine with gnu C/C++ compiler but have problems with Intel compilers.

Algorithm 6: JEIRA and GEISA pseudo code to determine the risk allele

Data: $caseMaxAllele$ is the most common allele in case and
 $controlMaxAllele$ is the most common in control
 $caseMaxRatio$ and $controlMaxRatio$ is the frequency of
 $caseMaxAllele$ and $controlMaxAllele$ respectively.

```
if  $caseMaxRatio > controlMaxRatio$  and
 $caseMaxAllele = controlMaxAllele$  then
     $riskAllele \leftarrow caseMaxAllele$ 
else
     $riskAllele \leftarrow caseMinAllele$ 
```

Boost 1.55 does not work with Intel compilers, 1.54 works. AllOf in Google Test also does not work together with Intel compilers.

Nvcc can not compile c+11 code, no interface included in the code that nvcc compiles can use c++11. This problem can largely be avoided by using separate compilation. The CUDA find package module for CMake(included in version 2.8 and later) has a default flag that should be removed as it can cause problems. The flag tells nvcc to propagate the host code flags, if the compiler for the host code then uses c++11 that flag is propagated to nvcc and it will not compile. The flag is disabled in CMake by using set(CUDA_PROPAGATE_HOST_FLAGS off).

Appendix D

How to Compile and Use CuEira

The source code for the program is available at github, <https://github.com/Berjiz/CuEira>.

List of dependencies. Listed version is the version used.

- Boost, tested with version 1.54, 1.55 does not work due to a bug with Intel compilers
- CMake, version 2.8
- CUDA, version 5.5
- MKL, version 13.1
- Intel compiler version 14.0
- Google Test and Google Mock, already included in the folder.

First install the dependencies, they might need to be compiled from source. Open a terminal and make a new directory where you want the program to be compiled. Enter the directory and write:

`cmake /path/to/CuEira/`

Then type:

`make`

The program should now be compiling. CMake might complain about not finding a dependency. If it does you need to tell CMake where it is. The path to boost is set by using `EXPORT BOOST_ROOT = /path/to/boost/`. The path to the compilers is done with `CXX = cmake /path/to/source` for the C++ compiler and `C = cmake /path/to/source` for the C compiler.